AD-A070 374     KANSAS UNIV   LAWRENCE DEPT OF COMPUTER SCIENCE          F/G 9/2
                IMPLEMENTATION AND EVALUATION OF INTERVAL ARITHMETIC SOFTWARE. --ETC(U)
                APR 79   J Q ARNOLD, F P FORD, R G HETHERINGTON   DACA39-76-M-0248
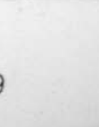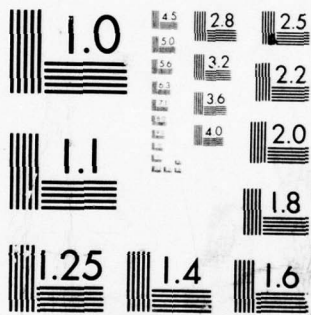                                        WES-TR-0-79-1                        NL

UNCLASSIFIED

1 OF 1

AD
A070374

END
DATE
FILMED
7--79
DDC

MICROCOPY RESOLUTION TEST CHART

NATIONAL BUREAU OF STANDARDS-1963-A

LEVEL

TECHNICAL REPORT O-79-1

# IMPLEMENTATION AND EVALUATION OF INTERVAL ARITHMETIC SOFTWARE

Report 3

## THE HONEYWELL G635 SYSTEM

by

James Q. Arnold, Frank P. Ford, Richard G. Hetherington

Department of Computer Science
University of Kansas
Lawrence, Kansas 66045

April 1979

Report 3 of a Series

DDC
JUN 26 1979
C

79 06 25 006

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER<br>Technical Report, O-79-1 | 2. GOVT ACCESSION NO. | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle)<br>IMPLEMENTATION AND EVALUATION OF INTERVAL ARITHMETIC SOFTWARE, Report 3. The Honeywell G635 System | | 5. TYPE OF REPORT & PERIOD COVERED<br>Report 3 of a series |
| | | 6. PERFORMING ORG. REPORT NUMBER |
| 7. AUTHOR(s)<br>James Q. Arnold<br>Frank P. Ford<br>Richard G. Hetherington | | 8. CONTRACT OR GRANT NUMBER(s)<br>Contract Nos.<br>DACA39-76-M-0248<br>DACA39-77-M-0107 |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br>University of Kansas<br>Department of Computer Science<br>Lawrence, Kans. 66045 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS<br>Integrated Software Research & Development Program, AT11 |
| 11. CONTROLLING OFFICE NAME AND ADDRESS<br>Office, Chief of Engineers, U. S. Army<br>Washington, D. C. 20314 | | 12. REPORT DATE<br>April 1979 |
| | | 13. NUMBER OF PAGES<br>88 |
| 14. MONITORING AGENCY NAME & ADDRESS(if different from Controlling Office)<br>U. S. Army Engineer Waterways Experiment Station<br>Automatic Data Processing Center<br>P. O. Box 631, Vicksburg, Miss. 39180 | | 15. SECURITY CLASS. (of this report)<br>Unclassified |
| | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE |

16. DISTRIBUTION STATEMENT (of this Report)

Approved for public release; distribution unlimited.

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)

18. SUPPLEMENTARY NOTES

19. KEY WORDS (Continue on reverse side if necessary and identify by block number)

Algorithms                          Honeywell G635 System
Computer systems programs           Interval arithmetic
Evaluation

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

This is Report 3 of a series entitled "Implementation and Evaluation of Interval Arithmetic Software." The series concerns implementation and evaluation of an interval arithmetic software package on six different computer systems. The other reports to be published in the series are:

Report 1: The State of the Interval: Evaluation and Recommendations

Report 2: The Honeywell MULTICS System

(Continued)

DD FORM 1473 EDITION OF 1 NOV 65 IS OBSOLETE
1 JAN 73

20.  ABSTRACT (Continued)

Report 4:  The IBM 370, DEC 10, and DEC PDP-11/70 Systems

Report 5:  The CDC CYBER 70 System

The most significant fact that has emerged from this project is that interval arithmetic has limited value as a tool for analyzing real algorithms.  The limitation is specifically dependency.

The magnitude and complexity of many of the problems being solved at the Waterways Experiment Station make interval analysis of the solution algorithms and sensitivity analysis of the data extremely difficult to accomplish with very high reliability using the current package.  The effort to employ interval arithmetic might better be directed toward development of new algorithms based on interval concepts than toward analysis of real algorithms currently in use.  Alternatively, redefining the arithmetic operations and interval representation might convert interval into a more practical tool for analysis of the real algorithms.

Accession For

NTIS  GRA&I
DDC TAB
Unannounced
Justification

By
Distribution/
Availability Codes

Dist   Avail and/or
       special

A

## PREFACE

In December 1975, the Automatic Data Processing (ADP) Center of the U. S. Army Engineer Waterways Experiment Station (WES), Vicksburg, Miss., submitted a proposal to implement and evaluate interval arithmetic, a software system for digital computer numerical analysis, on the Corps of Engineers' primary engineering computer—the WES Honeywell G635. The proposal was later expanded to include the implementation and evaluation of an interval arithmetic software package on six different computer systems. Engineering and scientific data problems were selected to be used on each of the six computers with the interval arithmetic software.

The work was funded by the Office, Chief of Engineers, U. S. Army, through the Integrated Software Research and Development (ISRAD) Program, AT11, Engineering Software Research.

This is Report 3 of a series entitled "Implementation and Evaluation of Interval Arithmetic Software." The other reports to be published in the series are:

Report 1: The State of the Interval: Evaluation and Recommendations

Report 2: The Honeywell MULTICS System

Report 4: The IBM 370, DEC 10, and DEC PDP-11/70 Systems

Report 5: The CDC CYBER 70 System

This report was written by Mr. James Q. Arnold, Mr. Frank P. Ford, and Dr. Richard G. Hetherington of the Department of Computer Science, University of Kansas, Lawrence, Kans. Their work was performed under Contract No. DACA39-76-M-0248, dated 28 April 1976, and Contract No. DACA39-77-M-0107, dated 24 February 1977, and was partially supported by the University of Kansas Computation Center project account and by Research Grant 3564-5038. The work concerned implementation and evaluation of an interval arithmetic software system on the Honeywell G635 computer system.

Dr. J. Michael Yohe, Director of Academic Computer Services, University of Wisconsin-Eau Claire, developed and wrote the interval arithmetic software package which was implemented on each of the six computer systems. Dr. Fred D. Crary, formerly with the U. S. Army Mathematics Research Center, University of Wisconsin-Madison, developed and wrote the AUGMENT precompiler which was implemented on each computer system as a front-end to the interval arithemtic software package. Dr. Yohe and Dr. Crary are specially thanked and recognized for their technical contributions and assistance.

Mr. James B. Cheek, Jr., formerly with the ADP Center, WES, provided initial impetus and guidance for the project. Mr. Fred T. Tracy, ADP Center, WES, provided expert advice and technical guidance during the project. Dr. N. Radhakrishnan, Special Technical Assistant, ADP Center, furnished technical guidance and general project supervision. The project and the report were monitored by Mr. William L. Boyt under the general supervision of Mr. D. L. Neumann, Chief of the ADP Center.

Directors of WES during the project and the preparation of the report were COL G. H. Hilt, CE, and COL J. L. Cannon, CE. Technical Director was Mr. F. R. Brown.

Copies of the other reports of the series, computer listings of the interval program and of AUGMENT for each computer system, and runs of the benchmarks for each computer system may be obtained from the ADP Center, WES.

## CONTENTS

3

## CONTENTS

4

IMPLEMENTATION OF THE 'AUGMENT' PRECOMPILER
AND 'INTERVAL' ON THE HONEYWELL G635

I.   PROJECT SUMMARY OF IMPLEMENTATION ACTIVITIES

AUGMENT version 4K was successfully implemented.

INTERVAL II was successfully implemented.

Test programs were run against the combined AUGMENT/INTERVAL II

   Package.

Some Problems were identified and some suggestions are made.

## II. INTRODUCTION

This report describes the Phase I activities of implementing AUGMENT and INTERVAL II on the Honeywell G635. As a result of the University's decision to release the G635 a month earlier than scheduled, Phase I had to be completed on the (compatible) Honeywell 66/60 which replaced the G635. Unfortunately, software errors on both Honeywell systems prevented completion of Phase I on schedule.

Part III of this report contains a discussion of the details of implementing AUGMENT, including a description of the host software errors mentioned above, observations and recommendations regarding efficiency, and arrangements for running AUGMENT. Part IV contains a discussion of the details of implementing INTERVAL II in the format suggested in Dr. Yohe's draft report Implementation of the Package on Other Hardware, June 9, 1976. Part V contains a description of the physical organization of the AUGMENT/INTERVAL II package and information related to use of the package.

Listings of the primitives written at K.U. are given in the Appendices, together with results of testing them. Also, the machine dependent constants needed by INTERVAL II are listed in an Appendix. Finally, a number of program listings together with their outputs are included in the Appendices.

All programs written at the University of Kansas and all modifications to the MRC programming packages to make them operational on the Honeywell systems have been carefully and thoroughly tested, and are believed to be correct. The University of Kansas and the programmers on this project disclaim responsibility for errors that may subsequently arise from use of the programs and systems described herein.

III.  IMPLEMENTATION OF AUGMENT

A.  DETAILS OF BRINGING AUGMENT UP

The initial effort to bring AUGMENT up was made on the University's
G635 which had 200K words (36 bits) for main memory, 380 million
characters of disc storage and both 7 and 9 track tape drives.  This
system supported batch and Time Sharing concurrently.  Wherever pos-
sible, the time Sharing System was employed in this project in order
to utilize on-line editing capabilities.  For files the size of AUGMENT
this approach was expensive and frequently ran into local boundary
limits on file space, processor time, and I/O,...the typical large file
problems.

Some minor difficulty was experienced in reading the original
tapes sent from MRC.  In spite of the obvious needs for standardization,
it seems still the case that vendor and installation prerogatives
combine to produce a 'tower of Babel' when it comes to exchanging tapes.
This was not a large problem, but is noted here for completeness and
to alert future users.

1.  The Primitives

After studying the documentation, it was decided to rewrite the
eight machine dependent primitives PACK, CCODE, MOVHOL, NUMIN,
ORDER, STRCHR, STRWDS and STRLNG.  These primitives had been
written for the G635 at the MRC and were not in the most efficient
form.  In particular, use of such functions as FLD was frequent,
and produced grossly inefficient object code on the G635.  All
eight primitives had been rewritten and were in the final stages
of debugging (only one known error still remained, albeit a
vexing one) when the G635 was released and the catastrophic system

7

errors on the 66/60 halted progress. Subsequently, the versions
of these eight primitives which had been sent from Wisconsin were
re-installed to eliminate all possible sources of programming
errors and this is the version currently being run. The more
efficient primitivies will be debugged as soon as possible and
installed as a part of the continuing tuning of the systems. A
copy of these routines will be made available to WES when they
have been thoroughly checked out.

2. MIN∅ and Variable Dimension (FDARGS)

Two relatively trivial errors appeared in the code sent from MRC.
MINO appeared where MIN∅ was required which was resolved by
having MINO call MIN∅. Also non-standard (for Honeywell, at least)
use of arguments in the routine FDARGS produced a fatal error.
Once identified, the correction was easy and was checked with
Dr. Crary.

3. Compiler Problems

During the implementation activity on both Honeywell systems,
three systems errors were identified. Two of the three are
probably exclusive with the new version (SR3I) on the 66/60 and
are not errors on version SR8F for the G635. The third error
resulted from a basic design flaw and exists on both the G635
and the 66/60 (and all other Honeywell 600, 6000, and 66/XX series
systems as far as is known).

   a. In outputting error messages, a list of subroutine calls
      in reverse order is provided by the FORTRAN compiler.
      Intermittantly, strings of zeros are placed in the name
      field of the list. In the midst of other chaotic errors,
      this caused some fruitless searches for array subscript

8

errors and investigations of table construction algo-
rithms in AUGMENT.  Apparently this error does not affect
the performance of AUGMENT, and should be ignored if
it occurs.

b.  A much more serious and disastrously time-consuming system
error occured in the computation of arguments to sub-
programs.  After two weeks of constant effort by the
entire project team, and aided by frequent long telephone
conversations with Dr. Crary in Wisconsin, it was dis-
covered that the 66/60 FORTRAN compiler was mishandling
the compilation of statements like

   PTR = SMMAKE(FSTCOL, ENDCOL, -IABS(CTYP), K, 0, -1, FALSE.)
(Statements of this form appear frequently in AUGMENT,
the one given being line AUG 38470 in PFETCH).  The
compiler (version SR2H on the 66/60) produced

   PTR = SMMAKE(-IABS(FSTCOL), K, 0, -1,.FALSE.)
Defining a dummy variable for -IABS(CØL) and assigning
it prior to compiling the error producing statement cir-
cumvented the problem.  This error should not appear
in SR8F on the G635.

c.  The third error concerns the assembly language instruction
Double Precision Floating Compare Magnitude (DFCMG).  If
DFCMG is used to compare two negative numbers whose mag-
nitudes differ by a large number (for example more than
$10^{25}$) the return is always "greater than".  In effect,
the argument in the register is always taken as having
magnitude greater than the argument in memory (even though,
for example, -1 might be in the register and $-10^{30}$ in
memory).  This error was first discovered on the G635 but

9

has subsequently been identified on all large scale Honeywell systems. Furthermore, the single precision form, FCMG, also fails as a result of the same basic design flaw. By Programming, these instructions were avoided. These errors are definitely part of the G635 system and should be publicized locally.

4. Testing the Package

The test programs sent with AUGMENT from the MRC were run successfully with all modifications indicated above.

One difficulty appeared in the way AUGMENT handles the following code:

```
INTERVAL A
        .
        .
        .
PRINT 100, A(1), A(2)
```

What happens is

(1) A is not taken over as INTERVAL in the AUGMENT output;

(2) the tables passed to the host FORTRAN are incorrect and produce the fatal error message that A has been previously defined as a scaler.

While it is admitted that the PRINT statement should be

```
PRINT  100, A ,
```

this is a potentially frequent error. It should produce an AUGMENT message and should not place incorrect information in the AUGMENT output.

10

B. LINKING

1. Size of the System

In dealing with the space-efficiency balancing problem the size
of the program AUGMENT may dictate the need for overlaying. In
the present implementation activity, no overlays were used. In-
stead the emphasis was on efficiency (especially run-time efficiency)
and large amounts of file and memory space were employed without
particular attention paid to the benefits of overlays. However,
in an ongoing-use environment, it is recommended that overlays
be employed to reduce costs and space requirements. This will be
necessary in a small memory installation; but is desirable in any
case. An overlay version is planned for implementation at K.U.
and will be available to WES.

2. Function Table Efficiency

The current version of AUGMENT generates function tables for each
use. Considerable efficiency might be gained by establishing
these function tables once and for all within AUGMENT and reducing
the overhead costs to each user.

C. PLAN FOR RUNNING AUGMENT

AUGMENT is a precompiler written in FORTRAN for FORTRAN
programs. It should be maintained as an object library on
tape or disk. It occupies approximately 600 blocks of disk
space. The library includes a mainline and two block common
initializing routines called BLDAT1 and BLDAT2. The control
cards needed to access the library are:

```
$            IDENT
$            ØPTIØN              FORTRAN, NØMAP
$            FØRTY
$            LIBRARY            LB
$            EXECUTE
$            FILE               LB
$            SYSØUT             21
$            FILE               20

*DESCRIBE

             (description deck)

*BEGIN

             (FORTRAN program)

*END

$            ENDJOB
```

This setup reads a Fortran program and a description deck
from file 05 (the usual input file) and outputs these on file 6.
File 21 is a copy of the output of AUGMENT which includes errors.
File 20 is an S* file which may be used as an input to the normal
Fortran compiler.

12

## IV. IMPLEMENTATION OF INTERVAL

Due to the later arrival from MRC of the INTERVAL tape, and at the suggestion of Dr. Yohe, attention was first given to the coding of the machine dependent arithmetic and bounding routines. After reading the documents on INTERVAL it was thought that the extended precision capability of the G635 could be exploited. Since all floating point arithmetic is done in the EAQ registers, one has available a full 72 bit mantissa which would allow guard digit computations. However, a little experience with the Honeywell rounding options was sufficient to demonstrate the pitfalls in that direction. For one thing, the floating store rounded instruction (FSTR) is incorrect in that it is possible to generate a number and its negative, which, after FSTR, do not sum to zero. FSTR always performs an "upward" round. A second problem results from the fact that the exponent has the same size for both single and double precision which precludes use of machine instructions in certain cases treated by INTERVAL. As a result, the machine rounding instructions were abandoned and the arithmetic routine BPAADD, BPASUB, BPAMUL and BPADIV as well as BPACEB (to convert EXTENDED to BPA) and BROUND (to handle directed roundings) were completely written in assembly language.

### A. DATA REPRESENTATION

Implicit assumptions that determine choice of data representation are listed by Dr. Yohe as

1. All operations and mathematical functions related to type BPA will be provided explicitly in the BPA package.

13

2. EXTENDED is used in evaluating the BPA mathematical functions, which requires binding EXTENDED to a higher precision data type than BPA.

3. A complete supporting package including all mathematical functions is required for type EXTENDED.

4. Every BPA number has an exact representation in type EXTENDED.

5. Every FORTRAN integer has an exact representation in type EXTENDED.

6. Conversion from REAL to BPA is exact.

7. Bounds on the accuracy of mathematical functions are available.

In the G635, the data type BPA has been taken as identical to REAL, and EXTENDED has been taken as identical to DOUBLE PRECISION. Although this greatly simplifies the data representation problems, items 1, 3 and 7 above are not readily met as a result of deficiencies in the Honeywell FORTRAN subroutine library and lack of usable documentation for what is there.

The following functions were not supplied in the vendors SR8F FORTRAN for the G635:

| BPA (REAL) | EXTENDED (DOUBLE PRECISION) | | |
|---|---|---|---|
| TAN | DTAN | DEXP2** | DTANH |
| CBRT | DARSIN | DSINH | DCBRT |
|  | DARCOS | DCOSH | DLOG2* |

\* Rename of DLOG

\*\* Rename of DEXP2

14

In the SR3I version of the 66/60 software, a new and complete set of FORTRAN subroutines are contained in the FORTRAN library. A copy of those routines was requested from Honeywell for delivery with this project. At the time of this writing it appears likely that the request will be granted. This may be handled at WES by installing the new library in place of the current one or by treating it as a special library for use with the AUGMENT/INTERVAL package. In any event, if the complete FORTRAN library is not used, the G635 FORTRAN library will not support INTERVAL unless locally developed subroutines are made available.

## B. CODING OF PRIMITIVES

This section contains a description of the algorithms developed for the BPA arithmetic and for BROUND, together with explanations of the decisions made. Initially these algorithms were written to incorporate special handling for the asymmetric numbers described below. However, after having obtained a working version which allowed the special cases to be treated correctly as far as INTERVAL was concerned, it was decided that these special cases caused inefficiencies and difficulties of understanding beyond their usefulness. For this reason a change was made to the algorithms which eliminates the special cases as either inputs to, or results from the arithmetic routines. The current versions contain these modifications as described below. However, without these cases, the requirement for passing arguments with the correct sign to BROUND is no longer justified. This in turn, implies that some improvement in efficiency will be realized by

15

a return to end sign correction as part of BROUND (this had been abandoned when it was observed that two passes through BROUND might be required for the special cases mentioned above) rewriting all the primitives in order to accomplsih this increase in efficiency is a possible step in tuning the system.

Normalized, floating point, two's-complement arithmetic on the Honeywell 635 dones not have symetric bounds on the numbers it can represent. For single precision, normalized numbers, the following ranges are given;

<u>negative</u>         $-(1+2^{-26})(2^{-129}) \geq N \geq -2^{129}$

<u>positive</u>         $2^{-129} \leq N \leq (1-2^{-27})(2^{127})$

Furthermore, zero is represented as $(0)(2^{-128})$.

Two main problems are caused by this lack of symmetry: the absolute value of the negative number with the largest magnitude cannot be represented, and the smallest, normalized, positive number has no normalized negative counterpart (although its negative can be given by an un-normalized number). To preserve the accuracy of the arithmetic done by INTERVAL, we must test for and identify these cases at various points in the computation.

Since the negative number with the largest magnitude does not have a negative, we do not allow it to be either an argument to or an answer from any of the arithmetic routines. If that number is passed to any of the arithmetic routines as an argument, the overflow indicator is set, the most negative number allowed by the arithmetic routines is loaded into the A register of Interval, and control is passed to BROUND to set the correct fault indicator (overflow or infinity) depending on the brounding option specified for that particular operation. If the negative

16

number with the largest magnitude is computed in any of the arithmetic routines, then it is treated like an overflow condition, which it actually is to Interval - even though it may not be an overflow to the computer.

Although there is a negative number with the same magnitude as the smallest positive number, it is not normalized, and therefore it cannot be returned as an answer from the arithmetic routines. Consequently, it is not allowed as an argument either. If the smallest, positive, normalized number comes in as an argument to the arithmetic routines, the exponent underflow indicator is turned on, and control passes to BROUND. In BROUND, it is treated as a case of underflow by one.

All arithmetic routines pass a number with the correct sign to BROUND. Corrections for the residue indicator also take place in the arithmetic routines.

> //In these arithmetic routines, MINNEG is the negative
> number with the largest magnitude, which is representable
> in the computer. NEGBND is the negative number with the
> largest number which is allowed in Interval. MINPOS is
> the smallest, normalized positive number which can be
> represented. Additional comments are given for numbered
> lines.//

BPASUB:    //compute  $A \leftarrow A1 - A2$ by negating A2 and adding //

   $U \leftarrow A2$

   $if$ $U = 0$

      $then$ $A \leftarrow A1$

         return

17

```
                    else U ← -U

                        goto SUBIN              //go to add routine to add A1,A2//

            fi

BPAADD:     //compute A ← A1 + A2//

            U ← A2

            if U = 0                            //if U = 0, then A1 is the answer//

                then A ← A1

                    return

            fi

SUBIN:      A ← A1

            if A = 0                            //if A = 0, then the answer is in U//

                then A ← U

                    return

            fi

1           if |U| < |A|                        //put argument with smaller//

                then  U ↔ A                     //magnitude in A//

            fi

            if A < 0                            //A must be made positive//

                then A ← -A

                    U ← -U

                    SC ← 1

            fi

            TEMP ← E(U) - E(A)                  //we compute the shift count//

2           if TEMP > 0                         //if TEMP ≤ 0, we need no shift//

3               then E(A) ← E(U)

4                   shift M(A) right TEMP bits          //line up mantissas//

            fi
```

```
          A ← A + U

          if SC = 1

              then if EO = 1

5                     then A ← -A

6                          EO ← 1

7                          EU ← 0

                      else if EU = 1

8                             then A ← -A

9                                  EU ← 1

10                                 EO ← 0

                              else A ← -A

                              fi

                      fi

          fi

          if RI = 0                          //if no residue and M(A) = 0//

              then if A = 0                  //then we don't need brounding//

                      then return

                      else goto BROUND

                  fi

          fi

          //now we need a residue correction//

11        if SC = 0

12            then if M(A)_{28-63} = 0

13                    then M(A)_{63} ← 1

                  fi

                  goto BROUND

14            else if M(A)_{28-63} = 0
```

19

15                              *then* M(A) ← M(A) − (0...01)

                    *fi*

                    *goto* BROUND

          *fi*                                          //end of BPAADD//


Additional comments:

<u>lines</u>

1        At the present time (1976) the DFCMG instruction does not

         work on the Honeywell 635; consequently, we must work

         around it, using other instructions to do the same thing. We

         have   done this by making both of the arguments positive,

         comparing them (which is the same as comparing their

         magnitudes), and then putting the smaller one in A and the

         bigger one in U with their correct signs.  The problem with

         the DFCMG was not known by Honeywell at the time it was

         reported to them.

2−4      E(U) − E(A) will be negative in one special case:

         M(A) = −M(U) = 010...0 = ½.  A negative number normalizes

         to have a 1 in the first bit of the mantissa; consequently,

         its exponent is one less than the exponent of its absolute

         value.  In this case, we do not want to shift the exponents.

         When the difference of the exponents is zero, even when one

         is negative and one is positive, we do not need to shift

         the mantissas, because they line up already.  Otherwise,

         we must shift the mantissa of the A register right, to

         line it up with the mantissa of the U register.  The

         number of bits we must shift it is the difference,

         E(U)−E(A).


20

5-10    After the addition of the two arguments, we must restore the

sign of the answer.  If we produced an overflow in the

addition, simple negation will produce the mantissa with

the correct sign.  One possible answer for the overflow

is the smallest positive normalized number.  If we negate

that, we will get an underflow, and so after the negation,

we must be sure to turn off the EU indicator, and to turn

the EO indicator back on.

Furthermore, if we produced an underflow in the addition,

one possible answer is the negative number with the largest

magnitude representable on the machine.  If we try to negate

that, we will get the correct mantissa, but we will also

produce an overflow because that number does not have an

absolute value which is machine representable.  Once

again, after we do the negation, we must be sure to turn

the EO indicator back off (if it was turned on), and we

must be sure to reset the EU indicator.

11-15   At this point we need to make a residue correction, and

we have four cases to deal with:

    1.  SC = 0 and A $\geq$ 0

    2.  SC = 0 and A < 0

    3.  SC = 1 and A $\geq$ 0

    4.  SC = 1 and A < 0

    Case 1.  Both A1 and A2 were positive.  Thus the bits

             lost decreased the magnitude of the answer.

             If  any of the bits $M(A)_{28-63}$ are 1, the

             correction has already been made, so only if

21

$M(A)_{28-63} = 0$ do we set $M(A)_{63} \leftarrow 1$.

Case 2. A1 and A2 had different signs, and the positive

had the smaller magnitude. Thus the bits

lost increased the magnitude of the answer.

If $M(A)_{28-63} = 0$, then we must set $M(A)_{63} \leftarrow 1$

to correct for the gain in magnitude. If

$M(A)_{28-63} \neq 0$, however, the correction has been

made.

Thus, Case 1 and Case 2 require identical treatment.

Case 3. A1 and A2 again differed in sign, but this time

the positive argument had the larger magnitude.

Therefore, the bits lost increased the

magnitude of the answer. If $M(A)_{28-63} \neq 0$

and A > 0, any truncation of non-zero bits in

$M(A)_{28-63}$ will insure proper brounding when

*augmentation is implied.* If $M(A)_{28-63} = 0$,

however, truncation will leave the answer too

large, and so we must subtract 0...01 from

$M(A)$ to insure proper brounding when truncation

is implied.

Case 4. Both A1 and A2 were negative, and thus the

non-zero bits lost decreased the magnitude

of the answer. If $M(A)_{28-63} \neq 0$, then

truncation will insure correct augmentation,

while the non-zero bits will guarantee proper

rounding towards zero. But once again,

if $M(A)_{28-63} = 0$, neither augmentation nor

22

brounding toward zero will work, and so we

must subtract 0...01 from M(A).

Thus Cases 3 and 4 require identical treatment.

BPAMUL: //compute A←A1*A2//

//convert arguments A1,A2 to positive values. Set SC:=1

if answer must be sign - corrected before bround//

*if* A1 < 0

*then* A1←-A1

SC←1

*fi*

*if* A2<0

*then* A2←-A2

SC←1-SC

*fi*

//compute A1*A2 without normalizing//

A←A1*A2

//handle overflow or underflow from above//

*if* exponent overflow *then go to* BPAMEØ *fi*

*if* exponent underflow *then go to* BPAMEU *fi*

A←normalized A

*if* exponent underflow from normalization *then go to* BPAMEU *fi*

*if* SC=1 *then* A←-A*fi*

*go* to BROUND

BPAMEU: A←-A

ØPTIØN←rounding option

BPAFLT ← 0

23

```
          if SC=1 then A←-A fi

          go to .EU.

          //.EU. is a section of BRØUND for exponent underflow//

BPAMEØ:   ØPTIØN ← rounding option

          BPAFLT ←0

          //normalize A.  This may fix exponent overflow//

             E(A)=0

          if SC=1 then A ←-A fi

             A←normalized A

          //E(A) has shift amount//

             E(A) = E(A) + E(A1) + E(A2)

          if exponent overflow then go to .EØ. fi

BPAD1V:   //compute A←A1/A2//

          //convert to positive values//

          if A1 < 0 then

             A1←-A1

             SC←1

          fi

          if A2=0 then go to ZERØDV fi

          if A2 <0 then

             A2←-A2

             SC←1-SC fi

          //determine if dividend larger than divisor.  This section

             (specially DVF instruction) works only if dividend ≤ divisor.//

          if M(A1)>M(A2) then

             M(A) ← M(A) /2

             SHIFT=1

          fi
```

24

$$\textit{if } M(A1) = M(A2)$$

$$SHIFT = 0$$

$$\textit{fi}$$

$A\leftarrow$ quotient $(M(A1)/M(A2))$

$Q\leftarrow$ remainder $(M(A1)/M(A2))$

$\textit{if } Q \neq 0 \textit{ then } RI:=1 \textit{ fi}$

//RI is residue indicator//

//negate if necessary.  Clear exponent to avoid its interference//

$E(A)=0$

$\textit{if } SC = 1 \textit{ then } A\leftarrow-A \textit{ fi}$

//now compute exponent. (A) is from normalization//

$E(A) = SHIFT$

$A \leftarrow$ normalized A

$E(A) = E(A) + E(A1) - E(A2)$

$\textit{if } E(A) > 128 \textit{ then } \textit{go to } EXP\emptyset \textit{ fi}$

$\textit{if } E(A) <-128 \textit{ then } \textit{go to } EXPU \textit{ fi}$

$\textit{if } RI = 1 \textit{ then } M(A)= M(A)+1 \textit{ fi}$

| | |
|---|---|
| EXP∅: | ∅PTI∅N ←0 |
| | BPAFLT ←0 |
| | *go to* .E∅. |
| EXPU: | ∅PTI∅N ← rounding option |
| | BPAFLT ←0 |
| | *go to* .EU. |
| ZERODV: | BPAFLT ←4 |
| | *go to* EXIT |

25

Understanding the brounding of two's complement numbers requires some familiarity with the way those numbers are represented, particularly, the negative numbers.

A positive, single-precision, normalized, floating-point number in the Honeywell 635 must have a 0 in bit zero, and a 1 in bit one. After that, it may have any combination of 1's and 0's. Consequently, the mantissa may represent a number as small as $2^{-1}$, and as large as $1-2^{-27}$. If we add bits to a positive mantissa or concatenate bits at the right end of a positive mantissa, and those bits are nonzero, then we increase the magnitude of the mantissa; if we drop a non-zero bit off the right end of a positive mantissa, we will decrease the magnitude of the number represented by the mantissa.

The Honeywell 635 has a 28-bit single-precision mantissa, and a 64-bit double-precision mantissa. The rounding options for Interval are accomplished as follows:

U,A: If any of the bits in $M(A)_{28-63}$ are non-zero, then we want to return the next higher single-precision mantissa as the answer. If $M(A)_{28-63} = 0$, though, we don't want to change the answer we have. Thus we want to add a number to A that will propagate a carry from any of the bits 28-63 into bit 27 of the mantissa without generating a carry if $M(A)_{28} = 0$.

L,T: Any non-zero bits in $M(A)_{28-63}$ make the answer too large, so we truncate them and use only the single precision answer. If the extra precision word is all zero, then we may simply truncate in that case too.

R: If we are rounding, a 1 in $M(A)_{28}$ means that we round

26

up, while a 0 means we round down. Thus, we may add a
one into $M(A)_{28}$ and we will get the correct answer.
If there is a one there, we will generate a carry
into the single precision word, if there is a zero,
then no carry will be generated, and the answer
will be rounded down as desired.

Negative numbers are somewhat different. First, they normalize
differently; they must have a 1 in the first bit, a zero in the second
bit, and after that, any combination of 1's and 0's may occur. Thus,
the range of negative mantissas is from $-1$ to $-(2^{-1} + 2^{-27})$. A negative
mantissa may represent a value with a greater magnitude than a positive
one, but a positive mantissa may represent a value with a smaller magnitude
than a negative one. If we add bits to a negative mantissa or concatenate
bits at the right end of a negative mantissa and those bits are non-zero,
we will decrease its magnitude; if we chop a non-zero bit off the
right end of a negative mantissa, we will increase the magnitude of the
number represented by the mantissa. This is very important to our
brounding. The brounding options are executed as follows:

U,T: Both U and T imply brounding toward zero; if any
of the bits in $M(A)_{28-63}$ are one, then we want to
generate a carry into $M(A)_{0-27}$ to decrease the
magnitude of the answer. Only if $M(A)_{28-63} = 0$
will we not want to generate that carry. Thus
we may add the same constant here as we did for
augmentation in the positive case.

L,A: To find the lower bound, any ones in $M(A)_{28-63}$ must
be ignored, because they decrease the magnitude of

27

the answer, and to find the lower bound (or to augment) we must increase the magnitude. If there are no 1's there, we have the correct answer. Thus for both L and A options here, we may simply chop off the double precision word.

R: If we are rounding a negative number, then we want to go away from zero if the extra precision word has a 1 in bit 28 and 0's in $M(A)_{29-63}$ or if bit 28 has a zero in it. To take the answer away from zero, we want to truncate those bits. On the other hand, if $M(A)_{28} = 1$, and any of $M(A)_{29-63} = 1$, then we want to bround toward zero. Consequently, we want to generate a carry into $M(A)_{27}$ of $M(A)_{28} = 1$ and any of the bits in $M(A)_{29-63} = 1$.

One important thing to realize about this brounding is that overflow and underflow may occur in certain conditions. We may get an overflow when augmenting positive numbers, and we may get an under-flow by one if we truncate negative numbers. These conditions must be checked when the brounding is done, and the appropriate steps must be taken if a fault is recognized.

Three numbers are used to do brounding: UTA is used for the U option (positive and negative), the T option (negative), and the A option (positive); NEGRND is used by the R option for negative numbers; POSRND is used by the R option for positive numbers. The mantissas of the three numbers are presented here.

28

```
UTA      [ 00...00 | 11...11 ]
          0      27  28     63

NEGRND   [ 00...00 | 011...11 ]
          0      27  28      63

POSRND   [ 00...00 | 100...00 ]
          0      27  28      63
```

When one of these numbers is used, its exponent is set to the exponent
of the number being brounded.

Finally, we may present the following decision table for the
brounding routine.

| OPTION | SIGN OF ANSWER | |
| | + | − |
| --- | --- | --- |
| U | add UTA | add UTA |
| L | none | none |
| T | none | add UTA |
| A | add UTA | none |
| R | add POSRND | add NEGRND |

In the following brounding routine, several variables are used,
and they represent certain floating-point constants.

MINNEG  is the most negative machine representable number.

NEGBND  is the most negative number allowed as an argument
        to or as an answer from the arithmetic routines.

MAXNEG  is the negative, single-precision, normalized
        number closest to zero.

MINPOS  is the smallest, normalized, positive number
        which can be represented in the machine.

POSBND  is the smallest, single precision, normalized
        number allowed as an argument or an answer from the

arithmetic routines.

MAXPOS   is the largest, positive number.

UTA, NEGRND, and POSRND are as described.

```
BROUND:  normalize A

         if CO ≠ 'U' and CO ≠ 'L' and CO ≠ 'T' and CO ≠ 'A' and CO ≠ 'R'

              then CO ← 'R'                    //if correction option is invalid, use R//

         fi

         case

          :EO = 1:

              if A ≥ 0

                   then A ← MAXPOS

                        if CO = 'A' or CO = 'U' or CO = 'R'

                             then INF ← 1

                        fi

                                  return

                   else A ← NEGBND

                        if CO = 'A' or CO = 'L' or CO = 'R'

                             then INF ← 1

                        fi

                                  return

                   fi

          :EU = 1:

              if  A ≥ 0

                   then if CO = 'A' or CO = 'U' or (CO = 'R' and underflow by 1)

                             then A ← POSBND

                             else A ← 0

                        fi

                        return
```

30

```
                        else if CO = 'A' or CO = 'L' or (CO = 'R' and underflow by 1)

                                then A ← MAXNEG

                                else A ← 0

                    fi

                    return

        fi

:CO = 'U':                          // find the upper bound on A//

    A ← A + UTA

    if EO = 1                       //overflow may occur when A is positive//

        then INF ← 1

             A ← MAXPOS

        else if EU = 1              //underflow may occur when A is negative//

                then A ← 0

             fi

    fi

    return

:CO = 'L':                          //find lower bound//

    return                          //store single-precision in all cases//

:CO = 'T':                          //bround toward zero//

    if A ≥ 0                        //store single-precision if positive//

        then return

        else A ← A + UTA    //A < 0//

             if EU = 1

                 then A ← 0

             fi

             return

    fi

:CO = 'R':                          //round to nearest machine number//

    if A ≥ 0
```

31

```
                    then A ← A + POSRND

                        if EO = 1

                            then A ← MAXPOS

                                INF ← 1

                    fi

                    return

            else A ← A + NEGRND                //A < 0//

                if EU = 1

                    then A ← MAXNEG

                fi

                return

        fi
    :CO = 'A':                                //augment A//

        if A ≥ 0

            then A ← A + UTA

                if EO = 1

                    then A ← MAXPOS

                        INF ← 1

                fi

                return

            else return                       //A < 0//

        fi

    end

//end of bround//
```

C.  MACHINE DEPENDENT CONSTANTS

The machine dependent constants, given in BPA format, which
must be provided to INTERVAL are listed in [1].  These constants
are located in COMMON/BPACON/ and for the G635 version the actual
numbers used are given in Appendix B.  While many of these re-
quired constants are obtained in a straightforward manner
(e.g. the left endpoint of the INTERVAL containing PI), several
must be derived from vendor documentation (e.g. half-length of
INTERVAL about 1 where LN accuracy decreases).  The approach
taken in determining these constants was the pragmatic one of
comparing output from the routines with published tables
[2,3,4,5,6].  It is expected that some improvement in the INTERVAL
results might be realized by spending more time on the evaluation
of these constants.

One additional comment in passing:  the description of the
constant FRACBD does not clearly bring out the fact that FRACBD is
that constant which when added to any number effectively truncates
the decimal part and produces a real integer.

D. ERROR BOUNDS FOR EXTENDED PRECISION ROUTINES

As described in [1], the method of implementing interval
mathematical functions is to assign an error term to each
function which is then utilized in the obvious way to determine
the value (interval) taken by the function. Let f be a con-
tinuous interval valued function of an interval variable which
we write as

$$f([a,b]) = [c,d].$$

Then, there exist points a' and b' in [a,b] such that $f(a') = c$
and $f(b') = d$. If the EXTENDED library routine (double precision
routine in the present instance) produces the value $f_E$ with an
error bound $\epsilon$, then the result computed by INTERVAL II is

$$\hat{f}([a,b]) = [\nabla(f_E (a')-\epsilon, \Delta(f_E (b')+\epsilon)]$$

where $\nabla$ and $\Delta$ are downward and upward directed roundings,
respectively.

From this formulation, it is seen that accurate error bounds
are essential to optimize interval function evaluations.
However, it is extremely difficult to calculate these bounds
accurately and hence the conservative approach suggested by
Dr. Yohe in [1] has been adopted. In effect this means using
the error bounds specified in the vendor documentation for these
routines.

34

E. MODIFICATIONS TO INTRAP

INTRAP provides a printed record of faults encountered
in interval operations. It is assumed that each routine
which can call INTRAP has at most three arguments, and
that the following information is provided to INTRAP:

ID  a 3-character suffix specifying the calling routine

TYPA  the type of argument A

TYPB  the type of argument B

TYPR  the type of the result

Values of TYPA, TYPB, TYPR may be any one of

    0 null

    1 FORTRAN INTEGER

    2 FORTRAN REAL

    3 FORTRAN DOUBLE PRECISION

    4 BPA

    5 EXTENDED

    6 INTERVAL

Because BPA was identified with REAL and EXTENDED was
identified with DOUBLE PRECISION, the built-in conversion
routines were suitable for this implementation and no
new code was added to INTRAP.


F. PROCESSING WITH AUGMENT

INTERVAL consists of a description deck for use with
AUGMENT and a package of compatible subroutines and functions
to handle INTERVAL data. It should be maintained as two
files--a library file and a description deck. The control

35

cards needed for it would be:

```
$    IDENT

     .
     .   (AUGMENT CONTROL CARDS)
     .

$    FILE      20, AISD
*  DESCRIBE
$    SELECT    INTERVAL DECK
*  BEGIN

     Fortran program

*  END
$    OPTION    FORTRAN, NOMAP
$    FORTY
$    LIBRARY   LC
$    FILE      S*,AIDD
$    EXECUTE
$    PRMFL     LC, INTERVAL LIB.
     Fortran data
$    ENDJOB
```

## G.  CHECKING THE PACKAGE

The test decks supplied by Dr. Yohe were run against
INTERVAL II in the manner described in the preceding section,
and the output was checked with the sample output.  As mentioned
in section E, these tests utilized INTRAP as it was sent
to us by MRC.  The machine dependent constants described in
C and D were used in these tests.  The results matched
and the test was assumed passed.  Appendix C contains these
test results.

H.  TUNING THE PACKAGE

Tuning the package is interpreted to mean making any modifications to code or data which will result in reducing time or space requirements, or will produce more accurate interval results.  (Conflicts may occur between these efforts; for example, run time may be increased in order to get tighter intervals.)  The major activities planned for tuning the system are:

1.  Replacing AUGMENT Primivitives discussed in section IIIA.1. of this report.

2.  Constructing an overlayed version of the package. (B-1 above)

3.  Fixing the AUGMENT function tables rather than regenerating them for each use.  (B-2 above)

4.  Improve error bounds for extended precision  routines (C and D above)

5.  Remove all unnecessary calls in the AUGMENT output (as suggested in [1])

6.  Rewrite arithmetic primitives for INTERVAL as discussed in III.B of this report.

V.   USING THE COMBINED AUGMENT/INTERVAL PACKAGE

A.   ORGANIZATION (PHYSICAL DESCRIPTION)

The particular organization to be used at WES will depend to great extent on local decisions.  The modules being delivered include

1.  Object library of AUGMENT programs.

2.  Source library of AUGMENT programs.

3.  Object library of INTERVAL programs.

4.  INTERVAL description deck for AUGMENT.

A perm file containing control cards (see III.C) needed to load and execute the AUGMENT object library and a perm file containing control cards (see IV.F) to run AUGMENT and the description deck will be created and accessed by $ SELECT cards.  Another perm file to describe the necessary libraries and the output file from AUGMENT to permit the FORTY compiler to execute the AUGMENT output will also be created and accessed by a $ SELECT card.


B.   USER INFORMATION

For input/output, it should be noted that a variable of type INTERVAL is converted to a variable of type REAL with dimension 2.  Therefore, normal FORTRAN I/O methods may be used.

Section IV F contains specific control card information for running programs.

EVALUATION OF INTERVAL ARITHMETIC
FOR THE HONEYWELL G635

## VI. INTRODUCTION

At its instigation this project was viewed as a relatively routine
implementation of the Augment/Interval programming system on the GE/
Honeywell 635 computer followed by a set of benchmark runs to investigate
the efficacy of using Interval Arithmetic to analyze error in the problem
solution methods for use at WES. From the beginning, difficulties were
encountered in all aspects of the project activity, including operating
system failures and misunderstood communications at K.U.; delays in
delivery of tapes, inefficiencies, and occassional errors in coding sent
by MRC; and untested benchmark programs and data sent by WES. While some
difficulties of this sort are to be expected, the sum total of all of
them coupled with the rather poor performance of the initial untuned
Interval/Augment system lead to extensive delays in the completion of
the contracted work. Every effort has been made to satisfy the letter
and the spirit of the contractual agreement on the parts of all parties
involved. However, the simple fact is that the task was monumentally
greater than anticipated and even now leaves many areas of potentially
fruitful further investigation unexplored.

Briefly, the chronology of events since the report on implementation
given November 1976 at Vicksburg is:

Corrected versions of the benchmark algorithms were
    received                                           Nov. - Dec. 1976

Overlay version of AUGMENT/INTERVAL implemented        January 1977

GAUSSE benchmark runs initiated after receiving
    the complete data                                  February 1977

SPLINE benchmark run successfully                      March 1977

| | |
|---|---|
| Report on completion of Phase III sent to WES | March 1977 |
| FFT benchmark runs initiated | April 1977 |
| INTERVAL arithmetic and brounding routines<br>    rewritten | May - June 1977 |
| BPA subroutine calls eliminated | July 1977 |
| Reruns of benchmarks for timing comparisons | July 1977 |

VII.   DISCUSSION OF DIFFICULTIES ENCOUNTERED IN USING INTERVAL ARITHMETIC

A.   DEPENDENCY

There is inherent in the operations of interval arithmetic a defect
which may increase width of intervals in computed results unrealistically.
The generic term describing this defect is dependency.  It derives from
endpoint calculations which make use of expressions which are dependent,
but the interval arithmetic operations treat them as independent.  Thus,
if $Y = \frac{X}{X}$ , where $X = [1/2,2]$, the definition of interval division produces
$Y = [\frac{1}{4},4]$.  This result includes the desired (correct) answer $[1,1]$ but
at the same time has increased the interval width substantially.  Other
anomolies include the fact that additive and multiplicative inverses
don't exist in the expected forms:  if $X = [a,b]$, $a \leq b$, then $X - X =$
$[a - b, b - a]$ which is not $[0,0]$ (unless $a = b$) and $X/X = [\frac{a}{b} , \frac{b}{a}] \neq [1,1]$
(unless $a = b \neq 0$).  The problem is that the two occurrences of the same
interval X, in these illustrations, are treated by the interval arith-
metic operations as independent when in fact they are dependent.  There
is no simple resolution to this difficulty because the definitions of
the interval-arithmetic operations must treat the two interval operands
as independent in order to guarantee inclusion of the correct result in
the interval answer.  Program testing of the operands for dependency is
hopelessly complicated in any but the simplest sequences of operations.
The fact is that interval arithmetic incorporates an error generation
characteristic which is distinct from ordinary arithmetic error considera-
tions.  The result is that proven methods and algorithms of real
arithmetic often fare badly when converted directly into interval
arithmetic.

41

The dependency problem is most effectively dealt with by avoidance –
simply not applying interval analysis to those algorithms or data which
are known to produce large dependency error. Unfortunately, there do not
yet exist practical methods for deciding, before the fact, what will be
the efficacious interval algorithms and what will not. Experimentation
often is the most direct means of obtaining such information although some
guidelines have been developed by F.N. Ris [13] and E.R. Hansen [12]. The
upshot of this situation for an installation such as WES is that many
production programs already in use cannot really be effectively analyzed
as interval arithmetic programs. In those instances where 'good' results
are obtained from the interval runs, one can be reasonably sure that the
real arithmetic algorithm is numerically stable. But when interval results
are meaningless because the intervals are too wide, then one cannot conclude
that the real arithmetic algorithm was also unstable – only that the increase
in interval widths was too great. If the dependency width (a term used by
F.N. Ris [13] to describe that part of the increased interval width due
exclusively to dependency) is the principal contributor to the increased
widths of interval results, then the algorithm should not be analyzed
using interval arithmetic and a different analysis technique should be
employed. On the other hand, if dependency width is small, even though the
interval width of results is too large, the algorithm may be salvaged by
some one or more of the traditional real arithmetic error reduction pro-
cedures, e.g., using higher precision calculations, more terms in the
approximations, etc. Of course, such modifications to an algorithm may
cause it to exhibit dependency which it previously did not have. To illustrate,
consider the simple series

42

$$f(x) = -1 + x - x^2 + x^3 - x^4 \ldots$$

If the first two terms are used, there is no dependency problem. But should two terms provide insufficient accuracy, adding the third term will improve accuracy and at the same time introduce a dependency problem.

B.   COMPARISON OF INTERVALS

Dependency is by far the overriding consideration in using interval arithmetic. If significant in an algorithm, it may obliterate any meaningful interval results. However, other sources of difficulty may enter the interval arithmetic operations for a particular solution method. One of these is the problem encountered in translating the compare operation from real arithmetic to interval arithmetic. An intuitive approach might be to compare endpoints and when the sup of one interval is less than the inf of the other interval, then say the first interval is less (smaller) than the second. This works fine for non-overlapping intervals. When the intervals overlap, an appeal might be made to a real comparison of their midpoints. But saying that the interval $[\frac{12}{25}, \frac{51}{100}]$ is less than the interval $[0,1]$, as would follow, leads to other difficulties since nearly half the numbers in $[0,1]$ are less than all values in $[\frac{12}{25}, \frac{51}{100}]$. The basic problem is that while real numbers are completely ordered, intervals are not. The result is that no entirely satisfactory way of translating a comparison of real numbers into a comparison of interval numbers exists.

C.   ZEROS ENTERING COMPUTED INTERVALS

Yet another source of difficulty in interval arithmetic stems from the fact that as a calculation progresses, the widths of the intermediate

interval results generally tend to become larger.  In particular, to guarantee the correct (real) answer is captured within the final interval, it is necessary to round intermediate answers to machine representable numbers which assure this inclusion.  While this strategy is obvious and correct, the effect is to expand interval width.  In most cases this somewhat conservative approach is defensible and does not lead to problems.  But when, in the course of a calculation, the intervals become null (contain the real number zero), or overlap, otherwise innocent operations may produce catastrophic interval growth.

D.    TIMING CONSIDERATIONS

Timing and space considerations can be viewed in a general way as well as particularizing them to a certain application program.  In this paragraph we will comment in general on these matters.  In the next section of this report more specific details of timing will be presented for each benchmark.  The efficiencies gained from tuning the package will be presented in the section following that.

Because the current version of AUGMENT/INTERVAL is designed for generality and portability, there are many inefficiencies in the package. Apart from the obvious increases in memory required for interval over real arithmetic, there are the time and space costs of generating function tables in AUGMENT for each use.  Modifying AUGMENT to eliminate function table generation with each use was investigated and found to be a substantial systems programming effort.  However, in an environment where there is frequent production use of AUGMENT, the effort to do the modifications would likely be justified.  This change has not been made in the current system.

44

Another inefficiency caused by the desire for generality is the multiple levels of subroutine calls involved in even the simplest source language statements. For example, the statement $X = A * B$, for A, B, X interval variables, generates a total of 134 subroutine calls, nested 3 deep.

Because in the H635 version BPA is identified with REAL and EXTENDED with DOUBLE PRECISION, we were able to reduce this sort of inefficiency by inserting replacement statements for subroutine calls. This work is described in Section IV of this report. Further effort along the lines of removing the nested subroutine calls could prove very beneficial to overall efficiency of the package.

Finally, some discussion with a faculty colleague has concerned the hardware implementation of interval arithmetic. There is no doubt that significant improvements in run times can be achieved by replacing the arithmetic primitives with hardware. However, dependency problems will not be rectified by hardware implementation.

VIII.   THE ENGINEERING ALGORITHMS

A.   SPLINE

A report made March 15, 1977 contains the results of the benchmark
run of the subroutine SPLINE.  These results are also included in this
final report for the sake of completeness.  As reported earlier, SPLINE
ran under INTERVAL with relative ease.  The interval widths of the
answers varied from 0.0003 to 0.0012 and the intervals were very nearly
centered on the real results.  Thus the SPLINE algorithm seems to be
fairly stable when performed in interval arithmetic and as remarked in
Section II of this report, this indicates stability of the algorithm
under real arithmetic.  The timing factor was 13.4, that is, it took
13.4 times as long to run under AUGMENT/INTERVAL as it did to run in
real arithmetic.  After tuning the system, as described in Section IV
of this report, the same run was repeated with the result that the
intervals were generally narrower and the timing factor was reduced to
12.2, a 10% improvement.  Listings of these runs are included in
Appendix B.

B.   FFT

Initial efforts to run the FFT benchmark failed as a result of an
error in the program modifications written at WES.  When this problem
was finally ferreted out in June, the FFT results were still not very
good.  Closer examination revealed an error in the INTERVAL cosine routine
sent to us by MRC.  After rewriting INTCOS we discovered that the same
fix had been distributed by MRC last summer and had been added to our
working system.  It is not clear how the correction was lost, nor whether

46

the uncorrected version was included in the package delivered to WES in November. However, the recently delivered system contains the corrected INTCOS.

Because of the difficulties mentioned above, a smaller problem had been constructed at K.U. to test FFT. It consists of an 8-point approximation to $\sin \frac{\pi X}{2}$ over the interval $[0,2]$. Interval results for this test appeared at first to have failed to capture the real answer, that is, the 'correct' answers were not in the intervals computed. After careful analysis of the algorithm and the programs, it was determined that the problem was in the real algorithm. The intervals had, in fact, captured the real arithmetic results. But the intervals were narrower than the error in the real number answers, and hence did not extend far enough to capture the 'true' answers (see results in Appendix B). In this way the interval analysis did in fact highlight the inaccuracies of the algorithm. However, these results also indicate that FFT is a numerically stable algorithm for this data. Going to double precision would probably improve accuracy of the real answers enough to have the intervals capture the 'correct' solutions. Unfortunately this conjecture cannot be tested since double precision is not available. But the conclusion regarding the stability of the FFT algorithm remains valid. Listings and output from the 8-point runs is included in Appendix B.

There was one significant change required in the benchmark driver. The program sent from WES included a calculation of the modulus of the answer. This calculation being the sum of squares of the real part (interval) and complex-part (interval) introduced a dependency error in the form of a negative argument for the square root function, INTSQT. To avoid this problem, the modulus calculation was taken out of the driver.

47

C.  GAUSSE

The first attempts to run GAUSSE with the data provided by WES all resulted in aborts due to timer runouts.  Even at a factor of 100, that is allowing the interval arithmetic run 100 times as much time as the real arithmetic run required, the program was completing only about two thirds of the problem.  Since this factor alone would deter using GAUSSE in an interval arithmetic setting, it was decided to investigate alternative means for analyzing the algorithm.  At this point we were unable to distinguish the causes of our problems, although it appeared that all of the difficulties described in Section II were contributing.  The data sent from WES was a randomly generated set of coefficients and right-hand constants for a 100 x 100 linear system, which we discovered later was incomplete (the last card had been lost from the deck).  When we reduced the problem to a 20 x 20 array of randomly generated input, the program ran to completion but produced useless results as reported in the Phase III report in March.

To gain better insight into the nature of the problems using INTERVAL with this algorithm, we decided to generate coefficients for smaller arrays which had properties known to produce 'good' results for real arithmetic GAUSS elimination with partial pivoting, and then to vary the data toward less well conditioned problems.  More than fifty runs were made using coefficient arrays generated as follows:

Upper left-hand element is called IONEONE

Diagonal elements increase by 1 going down the diagonal

First super diagonal contains -10

Second super diagonal contains -5

All other elements are -1

Solution vector is all +1

By varying IONEONE between 0 and 99900, we were able to study the behavior of the algorithm while controlling the condition of the problem. Several patterns emerged quite clearly.

(1) For well conditioned matrices, the interval widths are small and actually improve (become narrower) from the last to the first computed result. There is no pivoting. See run with IONEONE = 250 in Appendix B, for example.

(2) For less well conditioned matrices, pivoting usually takes place changing the order of solution, so it is not easy to follow the accuracy of successive solutions. The pattern that consistently emerges is that the intervals become larger starting with the last variable and progressing toward the first. However, the intervals decrease in width, starting five or six values before the variable for which pivoting takes place; and then, following the pivoted variable, increase exponentially (see example with IONEONE = 2, size = 40).

(3) For randomly generated data, the interval results are worse. For example, one 20 x 20 case produced intervals of virtually $(-\infty, \infty)$ for all variables, while the best random 20 x 20 case managed only 3 or 4 place accuracy.

Regarding timing, well conditioned systems ran under INTERVAL in about 16 - 18 times the real run time using the untuned system (reported in March). With the tuned version of the system, this factor was reduced to about 13. There was significant increase in time for INTERVAL runs using randomly generated data. For example, the 20 x 20 random case took 15 times the real case run time.

IX.   TUNING THE SYSTEM

A.   THE AUGMENT PRIMITIVES

After studying the documentation, it was decided to rewrite the
eight machine dependent primitives PACK, CCODE, MOVHOL, NUMIN, ORDER,
STRCHR, STRWDS and STRLNG.  These primitives had been written for the
G635 at the MRC and were not in the most efficient form.  In particular,
use of such functions as FLD was frequent, and produced grossly
inefficient object code on the G635.  All eight primitives were rewritten
in assembly language and installed in the package.  See Appendix C.

B.   LINKING

In dealing with the space-efficiency balancing problem, the size
of the program AUGMENT dictates the need for overlaying it to reduce
costs and space requirements.  This is essential in a small memory
installation but is desirable in any case.  The version of the system
delivered with this final report is appropriately linked for overlaying.

C.   PROVISION FOR MISSING FORTRAN LIBRARY SUBROUTINES

The AUGMENT/INTERVAL package requires the availability of a specific
set of subroutines in the host FORTRAN library.  Honeywell's FORTRAN
library for the 600 line does not include the following required routines
in a readily usable form:

| BPA (REAL) | EXTENDED (DOUBLE PRECISION) | | |
|------------|--------|--------|--------|
| TAN        | DTAN   | DLOG2* | DTANH  |
| CBRT       | DARSIN | DSINH  | DCERT  |
|            | CARCOS | DCOSH  | DEXP2**|

 *Rename of DLOG
**Rename of DEXP

In order to expedite implementation, dummy routines were inserted for these eleven functions. Subsequently, Honeywell released a new FORTRAN library on their 6600 line which included all of the above routines. However, the possibility exists that these newly released routines employ instructions from E.I.S. which are not available on the WES computer.

D.   MACHINE CONSTANT FRACBD

FRACBD is a machine dependent constant used in various INTERVAL routines. It is intended to provide a direct means for truncating the decimal part of a real number. The value assigned to FRACBD in the first implementation was close to but not exactly equal the number expected. As a result, there were occassional unexplained inaccuracies in the computations. The correct value has been inserted in the current version of the package. (See Appendix A).

E.   MODIFICATION OF INTMUL

In the subroutine INTMUL, as a result of some questionable coding, the variable TEMP could have been used in a logical comparison before it was assigned a value. This occurred in the statement

        IF(CASE .NE. 5 .OR. BPATMP(1) .LE. TEMP) GOTO 110

after the line labelled 104 and again in the statement

        IF(CASE .NE. 5 .OR. BPATMP(1) .GE. TEMP) GOTO 120

after the statement labelled 112. In both cases, if TEMP had not been assigned, then CASE was not equal to 5, and so the right side of the

51

condition did not have to be checked.  However, the Honeywell FORTRAN
compiler does execute the code which compares BPATMP(1) and TEMP.  Hence,
the contents of TEMP could cause errors.  By assigning zero to TEMP
upon entry to INTMUL, the possibility of error is eliminated.

F.    REMOVING UNNECESSARY CALLS TO BPA SUBROUTINES

Two changes were made to the description deck to increase the
efficiency of the INTERVAL package.

First, all calls to BPASTR were eliminated by removing the line

SERVICE COPY (STR)

from the BPA section of the deck.  The current version of INTERVAL
equates the types REAL and BPA, so A = B is used now instead of CALL
BPASTR(B,A).

Second, all calls to BPALT, BPALE, BPAEQ, BPANE, BPAGE, and BPAGT
have been eliminated.  Once again, since REAL and BPA are the same,
the output source from AUGMENT uses the corresponding FORTRAN logical
operators.

This was accomplished by adding an ENVIRONMENT section to the
description deck, and the FORTRAN logical operators were defined for
BPA operands.

To illustrate the effect of these changes, for the statement
X = A * B, noted in Section II D, the number of subroutine calls is
reduced from 134 to 46 and the nesting is only 2 deep.

52

G.   CODING OF THE PRIMITIVES

This section contains a description of the algorithms developed for BPA arithmetic and for brounding, together with explanations of the design decisions made.

Initially these algorithms were written to incorporate special handling for the asymmetric numbers described below.  However, after having obtained a working version which allowed the special cases to be treated correctly as far as INTERVAL was concerned, it was decided these special cases caused inefficiencies and difficulties of understanding beyond their usefulness.  For this reason a change was made to the algorithms which eliminated the special cases as either inputs to or results from the arithmetic routines (there is one exception, which will be explained later).  To insure accuracy during the brounding of the special cases, the original routine required the input to have the correct sign.  This requirement is no longer necessary, and so the current brounding routine requires non-negative input, and the sign correction is made at the end.

Normalized, floating-point, two's complement numbers on the Honeywell 635 are not symmetric about zero.  For single precision, normalized numbers the following ranges are given:

negative   $-(1 + 2^{-26})(2^{-129}) \geq N \geq -2^{129}$

positive   $2^{-129} \leq N \leq (1 - 2^{-27})(2^{127})$.

Furthermore, zero is represented as $(0)(2^{-128})$.

Two main problems are caused by this asymmetry:  the absolute value of the negative number with the largest magnitude cannot be

53

represented, and the smallest, normalized, positive number has no normalized complement (although it does have an unnormalized complement). As mentioned above, these two numbers are not allowed as inputs to the arithmetic routines. Thus every valid argument has an exact complement in the INTERVAL package.

No arithmetic operation takes place if either of these two numbers is found. When the most negative number is an argument to the routines, the infinity or the overflow indicator is set, and the negative number with the largest magnitude allowed is returned as the answer. When the smallest positive number is used as an argument, the underflow indicator is set, and either zero or the smallest positive number that is allowed is returned as the answer, depending on the browsing option. The single exception to this treatment occurs in division by zero, in which the unchanged dividend is returned as the "answer," and the zero divisor indicator is set.

Understanding the following routines requires some knowledge of how the floating-point numbers are represented.

A positive, single precision, normalized, floating-point number has an 8-bit exponent and a 28-bit mantissa. The exponent part holds a two's complement integer, E, and the mantissa part holds a fractional, two's complement number, M. The relation for a floating-point number Z is this:

$$Z = M * 2^E.$$

Furthermore, the mantissa must have a zero in the first bit (the sign bit), and a 1 in the second bit. Thus if M is positive, then $2^{-1} \leq M \leq 1 - 2^{-27}$. If we catenate non-zero bits to the right end

54

of a positive mantissa, then the magnitude of the number represented is increased; if we truncate non-zero bits from the right end of a positive mantissa, then the magnitude is decreased.

Double precision numbers also have 8-bit exponents, but their mantissas are expanded to 64 bits. The arithmetic and brounding routines use the full AQ register to manipulate the $M(A)$, and so a single precision mantissa occupies $M(A)_{0-27}$ with $M(A)_{28-72}$ being initially set to zero. The U register is simulated by a word pair in memory.

Given the elimination of the normalized numbers without (normalized) complements, the range of numbers allowable in the INTERVAL package is the following:

<u>negative</u> $-(1 + 2^{-26})(2^{-129}) \geq N \geq -(1 - 2^{-27})(2^{127})$

<u>positive</u> $(1 + 2^{-26})(2^{-129}) \leq N \leq (1 - 2^{-27})(2^{127})$.

Four machine dependent constants are used in the following routines.

MINNEG    is the negative number with the largest magnitude.

MINPOS    is the smallest positive number.

POSBND    is the smallest positive number allowed in the INTERVAL package.

MAXPOS    is the largest, single precision, positive number.

In addition, two addition constants are used by the brounding routine, and their mantissas are given below.

| ROUND | 00 . . . . . 00 | 100 . . . . . 00 |
|-------|-----------------|------------------|
|       | 0            27 | 28            63 |

| AUGMNT | 00 . . . . . 01 | 00 . . . . . 00 |
|--------|-----------------|-----------------|
|        | 0            27 | 28           63 |

When one of these addition constants is used, its exponent is set to
the exponent of the number to which it is being added.  The mantissas
of ROUND and AUGMNT are not normalized.

The arithmetic routines BPACEB, BPASUB, BPAADD, BPAMUL, and
BPADIV, and the brounding routine BROUND are given below.  Additional
comments are given for the numbered lines.

```
    BPACEB :   <<Convert the double precision Al to single
               <<precision and bround it.
               <<RES ← Al                              >>
        clear all indicators ;
        A ← Al ;
        case
            : A = MINNEG :  <<Overflow.>>
                 SC ← 1 ; goto .EO. ;
            : A = MINPOS :  <<Underflow.>>
                 goto .EU. ;
        endcase ;
        if A < 0
            then A ← -A ;  SC ← 1 ;
        fi ;
        goto BROUND ;  <<End of BPACEB>>


    BPASUB :   <<Compute RES ← Al - A2.>>

        clear all indicators ;
        U ← A2 ;
        if U = 0
            then <<The answer is Al.>>
              A ← Al ;
              case
                  : A = MINNEG :  <<Overflow.>>
                       SC ← 1 ; goto .EO. ;
                  : A = MINPOS :   <<Underflow.>>
                       goto .EU. ;
              endcase ;
              RES ← A ; return ;

        fi ;

        case
            : U = MINNEG :  <<Overflow.>>
                  SC ← 1 ; goto .EO. ;
            : U = MINPOS :   <<Underflow.>>
                  A ← U ; goto .EU. ;
        endcase ;
```

56

```
        U ← -U ;
        goto SUBIN ; <<BPAADD completes the subtraction.>>

        BPAADD :   <<Compute RES ← A1 + A2.>>

             clear all indicators ;
             U ← A2 ;
             if U = 0
                   then <<A1 is the answer.>>
                   A ← A1 ;
                   case
                        : A = MINNEG :   <<Overflow.>>
                               SC ← 1 ; goto .EO. ;
                        : A = MINPOS :   <<Underflow.>>
                               goto .EU. ;
                   endcase ;
                   RES ← A ; return ;

        fi ;

        case

             : U = MINNEG :   <<Overflow.>>
                   SC ← 1 ; goto .EO. ;
             : U = MINPOS :   <<Underflow.>>
                   A ← U ; goto .EU. ;
        endcase ;

        SUBIN :   <<This is the entry point for the subtraction routine.   >>
             A ← A1 ;
             if A = 0
                   then <<U holds the exact answer.>>
                        RES ← U ; return ;

        fi ;

        case

             : A = MINNEG :   <<Overflow.>>
                   SC ← 1 ; goto .EO. ;
             : A = MINPOS :   <<Underflow.>>
                   goto .EU. ;
        endcase ;

        if |U| < |A|
             then U ↔ A ;
        fi ;
        if A < 0
             then <<Make A positive.>>
                   A ← -A ; U ← -U ; SC ← 1 ;
        fi ;
        TEMP ← E(U) - E(A) ; <<Compute shift count.>>
```

```
        if TEMP > 0 <<No shift if TEMP < 0.>>
            then
                E(A) ← E(U) ;
                shift M(A) right TEMP bits ; <<Line up the mantissas.>>
        fi ;

        A ← A + U ;
        if A = 0
            then <<The answer must be exact.>>
                RES ← A ; return ;
        fi ;
        if RI = 1

            then <<We lost digits in the mantissa, and the
                  <<number in A is too small, regardless of
                  <<whether it is positive or negative.  We
                  <<set bit 35 to 1; this bit is in the extra
                  <<precision word, and will not affect the
                  <<rounding option.                          >>

            M(A)₃₅ ← 1 ;

        fi ;

        if A < 0

            then A ← -A ; SC ← -SC ;

        fi ;

        goto BROUND ; <<End of BPAADD>>
```

Additional comments for BPAADD:

Lines

1  At the present time (1976-1977) the floating point compare

   magnitude instruction DFCMG does not work on the Honeywell

   635; consequently, we program around it by making both

   arguments positive, and then comparing them.  The problem

   with DFCMG was not known by Honeywell at the time it was

   reported to them.

2-4  The exponent of the absolute value of an exact power of two

   is always one greater than the exponent of its complement;

consequently, $E(U) - E(A)$ will be negative if we are adding
a power of two and its complement. If the difference in the
exponents is negative or zero, then no shifting need be done.
Otherwise, we need to live up the mantissas for the addition and
change exponent of A accordingly. The residue indicator may
be set during the shift.

5      If the computed result is MINPOS, then this negation will cause
the exponent underflow indicator of the machine to be set,
and then it is treated as an underflow-by-one in BROUND. The
result can not be MINNEG.

```
BPAMUL :   <<Compute RES ← A1 * A2.>>

    clear all indicators ;

    U ← A2 ;

    if U = 0 or A1 = 0

        then RES ← 0 ; return ;

    fi ;

    case

          : U = MINNEG : <<Overflow.>>
                SC ← 1 ; goto .EO. ;

          : U = MINPOS : <<Underflow.>>
                A ← U ; goto .EU. ;

    endcase ;

    if U < 0
        then <<Make it positive.>>
                U ← -U ; SC ← 1 ;
    fi ;

    A ← A1 ;

    case
```

```
            : A = MINNEG :   <<Overflow.>>
                  SC ← 1 ; goto .EO. ;

            : A = MINPOS :   <<Underflow.>>
                  SC ← 0 ; goto .EU. ;

         endcase ;

         if A < 0

            then <<Make it positive.>>

                  A ← -A ; SC ← -SC ;

         fi ;

1        A ← A * U ;

         goto BROUND ; <<End of BPAMUL>>
```

Additional comments for BPAMUL:

Lines

1   Single precision mantissas have 28 bits (27 bits plus a sign),

    and so the product can be held in 55 bits ( a sign bit plus

    54 for the mantissa).  FMP computes the full AQ (72 bits)

    accurately, and so we use the machine's floating multiply

    instruction.

```
BPADIV : <<Compute RES ← A1/A2.>>

         clear all indicators ;

         U ← A2 ;

         if U = 0

            then <<Division by zero.>>

                  BPAFLT ← 4 ; <<DZ ← 1>>

                  RES ← A1 ; <<No checking here.>>

                  return ;
```

```
                    fi ;

                    if A1 = 0

                         then RES ← 0 ; return ;

                    fi ;

                    case

                         : U = MINNEG : <<Overflow.>>
                              SC ← 1 ; goto .EO. ;

                         : U = MINPOS : <<Underflow.>>
                              A ← U ; goto .EU. ;

                    endcase ;

                    if U < 0

                         then <<Make it positive.>>

                              U ← -U ; SC ← 1 ;

                    fi ;

                    A ← A1 ;

                    case

                         : A = MINNEG :  <<Overflow.>>
                              SC ← 1 ; goto .EO. ;

                         : A = MINPOS ;  <<Underflow.>>
                              SC ← 0 ; goto .EU. ;

                    endcase ;

                    if A < 0

                         then <<Make it positive.>>

                              A ← -A ; SC ← -SC ;

                    fi ;

                    V ← 0 ;

          for j ← 35 downto 1 do

  1       if |M(A)| ≥ |M(U)|

                    then
```

61

```
2                          M(A) ← M(A) - M(U) ;
3                          V ← V + 1 ;

            fi ;

4           shift M(A) and V left by 1 ; <<Multiply by 2.>>

       endfor ;

       shift M(A) right 36 bits ; <<Signal remainder.>>

       M(A) ← V ; <<Move the quotient to M(A).>>

       if M(A)_0 = 1

            then <<We need to normalize the quotient.>>

                  shift M(A) right 1 bit ;

                  E(A) ← E(A) + 1 ; <<Correct the exponent.>>

       fi ;

5           E(A) ← E(A) - E(U) ;

6           if overflow

                  then <<We may have exponent overflow or underflow.    >>

7                       if E(A)_0 = 1

                             then <<Overflow.>>

8                                  goto .EO. ;

                             else <<Underflow.>>

9                                  goto .EU. ;

                        fi ;

            fi ;

            if RI = 1

                  then <<Set a low order bit in M(A) to 1.              >>

                        M(A)_{35} ← 1 ;

            fi ;

            goto BROUND ; <<End of BPADIV.>>
```

62

Additional comments for BPADIV:

<u>Lines</u>

1-4     Yohe's algorithm for division in the base $\beta$ number system
        includes an extra loop at this point to subtract M(U)
        from M(A) until M(A) < M(U).  The normalized form of
        floating-point numbers in the Honeywell 635 makes this
        extra loop unnecessary; the M(A) must always be less than
        2 * M(U), and so the subtraction in line 2 insures that
        M(A) < M(U).

5-6     If the subtraction in line 5 produces a number either too
        large or too small to be represented in the 8-bit exponent
        of the 635, then the overflow indicator of the machine is
        set.

7-9     To determine whether the fault was caused by exponent
        overflow or underflow, we inspect the exponent computed.
        If E(A) is too big, then the machine uses the sign bit to
        gain the magnitude it needs.  Thus exponent overflow sets
        $E(A)_0$ to 1.  If E(A) is too small, then $E(A)_0$ becomes zero
        when the subtraction is performed.


BROUND : <<Bround the number in A, and return the answer in RES.
         <<The brounding option is held in OPTION.  The fault
         <<value is returned in BPAFLT, and the following
         <<values are used:

         <<            0    operation successful
         <<            1    exponent overflow
         <<            2    infinity
         <<            3    exponent underflow
         <<            4    zero divisor                        >>

```
if OPTION ∉ {'U', 'L', 'T', 'R', 'A'}

    then <<Use R if OPTION is invalid.>>

          OPTION ← 'R' ;

fi ;

if exponent overflow or A > MAXPOS

    then goto .EO. ;

fi ;

if exponent underflow or A < POSBND

    then goto .EU. ;

fi ;

      <<The correction for the residue indicator has been made
      <<already, and is held in M(A)₃₅.                        >>

case

1     : [(OPTION = 'U' and SC = 0) or

2      (OPTION = 'L' and SC = 1) or

3      OPTION = 'A'] and M(A)₂₈₋₆₃ ≠ 0 :

          RES ← A + AUGMNT ;

          if exponent overflow

              then <<Set the infinity flag and load MAXPOS.>>

                    BPAFLT ← 2 ;

                    RES ← MAXPOS ;

          fi ;

          if SC = 1

              then RES ← ⌐RES ;

          fi ;

          return ;
```

```
4       : (OPTION = 'U' and SC = 1) or

5         (OPTION = 'L' and SC = 0) or

6          OPTION = 'T' or M(A)_{28-63} = 0 :

                RES ← A ;

                if SC = 1

                    then RES ← -RES ;

                fi ;

                return

        : OPTION = 'R' :  <<the last possibility>>

                RES ← A + ROUND ;

                if exponent overflow

                    then <<set the infinity flag and load MAXPOS>>

                            BPAFLT ← 2 ;

                            RES ← MAXPOS ;

                fi ;

                if SC = 1

                    then RES ← -RES ;

                fi ;

                return ;

    end case ;

        .EO. :  <<Overflow faults are handled here.  The input
                <<number is irrelevant.                          >>

        case

            : (OPTION = 'U' and SC = 1) or

              (OPTION = 'L' and SC = 0) or

               OPTION = 'T' :
```

```
                    <<The brounding option implied truncation,
                    <<and so use the fault value for exponent
                    <<overflow.                                    >>

                    BPAFLT ← 1 ;

              else <<Augmentation was implied, and so we had
                    <<an infinity fault.                           >>

                    BPAFLT ← 2 ;

        endcase ;

        RES ← MAXPOS ;

        if SC = 1

            then RES ← -RES ;

        fi ;

        return ;

  .EU. :   <<Underflow faults are handled here.  Special
           <<action is taken if the R option is specified
           <<and we have exponent underflow by one.             >>

        BPAFLT ← 3 ; <<exponent underflow>>

        case

              : (OPTION = 'U' and SC = 0) or

                (OPTION = 'L' and SC = 1) or

                 OPTION = 'A' :

                        << The brounding option implied
                        << augmentation, and so we use the
                        << non-zero number with the smallest
                        << magnitude allowed.                     >>

                        RES ← POSBND ;

              : OPTION ≠ 'R' :

                        <<Truncation was implied; so we
                        <<use zero.                               >>

                        RES ← 0 ;
```

66

```
        else <<OPTION = 'R'>>

7           if A ← POSBND or exponent underflow by 1

                then RES ← POSBND ;

                else RES ← 0 ;

        fi ;

    endcase ;

    if SC = 1

        then RES ← -RES ;

    fi ;

    return ; <<End of BROUND>>
```

Additional Comments for BROUND:

Lines

1-3          Augmentation is implied in 3 cases:  (1) OPTION = 'U'
             and SC = 0 means the true answer is positive, and we
             want the least upper bound; (2) OPTION = 'L' and SC = 1
             means the true answer is negative and we want the
             greatest lower bound, thus we must augment the positive
             number to increase the magnitude; (3) OPTION = 'A'.
             If $M(A)_{28-63} = 0$ then the result is exact.

4-6          Truncation is implied in 4 cases:  (1) OPTION = 'U'
             and SC = 1 means the true result is negative, and we
             want the least upper bound, thus we truncate the positive
             result to decrease the magnitude; (2) OPTION = 'L' and
             SC = 0 means the true result is positive, and we want
             the greatest lower bound; (3) OPTION = 'T'; (4) $M(A)_{28-63} = 0$

67

4-6        means the result can be expressed exactly in single

precision, and no brounding option will affect it,

thus the extra precision bits are effectively truncated.

7        Exponent underflow-by-one can occur two ways: by

computation and by input of MINPOS. All numbers

smaller than POSBND and greater than or equal to MINPOS

are considered to be cases of underflow by one. If we

actually cause a machine exponent underflow during

computation, the computed exponent is 127, or the

largest exponent possible.

X. DESCRIPTION OF THE TAPE SENT TO WES

The tape sent to WES is a multifile tape, written at 800 bpi, 7 track.
The fifteen files contained on the tape are:

File 1    Arithmetic and brounding routines for INTERVAL.

File 2    Machine dependent primitives for AUGMENT

File 3    SESOL and BANSOL

File 4    INTERVAL source library, including corrected INTCOS and the
          modifications to the AUGMENT description deck to remove the
          superfluous subroutine calls

File 5    FFT

File 6    SPLINE(INTERVAL form)

File 7    GAUSSINT (old version of GAUSE)

File 8    GAUSS (INTERVAL form)

File 9    SESOL

File 10   FFT (original)

File 11   FFT (no complex)

File 12   FFT (INTERVAL)

File 13   FFT (8 point real)

File 14   FFT (8 point INTERVAL)

File 15   GAUSS (a third version)

XI. CONCLUSIONS AND RECOMMENDATIONS

In my opinion, the most significant single fact that has emerged
from this project is that interval arithmetic has limited value as a
tool for analyzing real algorithms.  The limitation is specifically
dependency.  As pointed out by our results and, in a more general way,
by Ris [2], the fact that interval arithmetic gives intervals which
are too wide does not tell you, necessarily, that the real algorithm
is bad nor that the real algorithm is necessarily sensitive to the real
data used.  In particular, it is well known that Gauss elimination with
partial pivoting is a stable solution technique in a real number setting
provided one is careful about conditioning and error control.  What is
not so well known is that the same Gauss method is subject to severe
dependency problems which can cause as much as fourfold increases in
interval widths at each stage of the procedure when employing interval
arithmetic.  The point is that interval analysis yields information about
interval algorithms in all cases; but interval analysis yields information
about the real version of an interval algorithm only when the interval
algorithm is 'good'.  One needs a way to break down the sources of
increased width when the computed intervals are too wide in order to be
able to say anything about the corresponding real algorithm.  Even then,
the conclusion may be that nothing can be inferred from the interval
analysis.  Both Ris [2] and Hansen [1] attack this problem, but from
different points of view.  Hansen proposes a generalized interval arith-
metic (g.i.a.) to reduce the effect of dependency.  By standardizing
the form of interval representation and redefining the arithmetic
operators, he is able to reduce the inherent lack of sharpness (caused
by dependency) to a second order effect.  Ris tackles the problem in a

70

different way, by defining predicates and functions on intervals which have predictable propagation properties in interval arithmetic operations. He then applies these ideas to analyze interval algorithms, highlighting those characteristics which might be used to predict the success or failure of the algorithm. Hansen's approach is more in keeping with the current project in that it reduces the factor (dependency) that limits the value of interval analysis of real algorithms. Ris' approach is designed to develop new algorithms for which interval arithmetic is well suited. Of course, once a 'good' interval algorithm is produced, it will be a 'good' real algorithm.

Getting more specific, the results of the benchmark runs indicate the following:

SPLINE is reasonably stable for the data run

FFT is very stable but suffers from loss of accuracy in the calculations for the 4096 point case.

GAUSSE is (predictably) stable for well conditioned matrices and becomes less so as the condition of the matrix deteriorates. For randomly generated data, the algorithm is less efficient and less accurate. Pivoting greatly increases widths due to dependency.

Recommendations:

In view of the number and range of difficulties encountered in attempting to use this AUGMENT/INTERVAL package for analyzing algorithms, it clearly should be treated as an experimental tool, and not considered for routine or production use in its present form.

The magnitude and complexity of many of the problems being solved at WES makes interval analysis of the solution algorithms and sensitivity analysis of the data extremely difficult to accomplish with very high reliability using the current package. A major effort to employ interval

71

arithmetic might better be directed toward development of new algorithms based on interval concepts (as implied in Ris' work) than to attempt analysis of the real algorithms currently in use. Alternatively, redefining the arithmetic operations and interval representation (as Hansen suggests) might convert INTERVAL into a more practical tool for analysis of the real algorithms.

72

# REFERENCES

[1]  Yohe, J. M.  "Guide to Implementation of the Interval II Package on Other Hardware," Draft report, Jun 1976, Mathematics Research Center, The University of Wisconsin-Madison.

[2]  Abramovitz, M. and Stegun, I. A., Handbook of Mathematical Functions, AMS 55, National Bureau of Standards, 1964.

[3]  National Bureau of Standards, Tables of the Exponential Function, 1939.

[4]  National Bureau of Standards, Tables of Natural Logarithms, Vol. III, 1941.

[5]  National Bureau of Standards, Tables of Circular and Hyperbolic Sines and Cosines, 1940.

[6]  National Bureau of Standards, Tables of Circular and Hyperbolic Tangents and Cotangents, 1947.

[7]  Crary, F. D., "Language Extensions and Precompilers," Technical Summary Report No. 1317, 1973, Mathematics Research Center, The University of Wisconsin-Madison.

[8]  Yohe, J. M., "Best Possible Floating Point Arithmetic," Technical Summary Report No. 1054, 1970, Mathematics Research Center, The University of Wisconsin-Madison.

[9]  Crary, F. D., "The AUGMENT Precompiler; I:  User Information," Technical Summary Report No. 1469, 1974 (revised 1976), Mathematics Research Center, The University of Wisconsin-Madison.

[10]  Crary, F. D., "The AUGMENT Precompiler; II:  Technical Documentation," Technical Summary Report No. 1470, 1975, Mathematics Research Center, The University of Wisconsin-Madison.

[11]  Ladner, T. D. and Yohe, J. M., "An Interval Arithmetic Package for the UNIVAC 1108," Technical Summary Report No. 1055, 1970, Mathematics Research Center, The University of Wisconsin-Madison.

[12]  Hansen, E. R., "A Generalized Interval Arithmetic," Interval Mathematics, No. 29 in the series Lecture Notes in Computer Science, Springer-Verlag, 1975.

[13]  Ris, F. N., "Tools for the Analysis of Interval Arithmetic," Interval Mathematics, No. 29 in the series Lecture Notes in Computer Science, Springer-Verlag, 1975.

The following are additional computer listings and runs available from the Automatic Data Processing Center, U. S. Army Engineer Waterways Experiment Station, P. O. Box 631, Vicksburg, Miss. 39180:

Interval Primitives written at the University of Kansas
Augment Primitives written at the University of Kansas

Test programs and their results:

Summation of first 129 terms of $(1/X)**(I-1)$
Roots of a quadratic equation
635 rounding errors in addition
Test mathematical functions using INTRAP
Gaussian elimination routine-interval extension
Table of factorials and their natural logarithms
SPLINE (real and interval)
FFT (8-point, real arc interval)
FFT (4096-point, real and interval)
Gaussian (20X20, real and interval)

APPENDIX A:  MACHINE DEPENDENT

CONSTANTS FOR 'INTERVAL II'

The following constants, in BPA format, have been provided and are all located in COMMON /BPACON/

| NAME | VALUE | DESCRIPTION |
|---|---|---|
| PIO2I | Ø002622077325 | Left endpoint of interval containing PI/2 |
| PIO2 | Ø002622077326 | Right endpoint of same |
| PII | Ø004622077325 | Left endpoint of interval containing PI |
| PI | Ø004622077326 | Right endpoint of same |
| TPII | Ø006622077325 | Left endpoint of interval containing 2*PI |
| TPI | Ø006622077326 | Right endpoint of same |
| THPI | Ø010455457440 | Right endpoint of interval containing 3*PI |
| ZRO | 0.0 | zero |
| ONE | 1.0 | 1 |
| ONM | -1.0 | -1 |
| TWO | 2.0 | 2 |
| BPAMNB | Ø400400000000 | smallest positive BPA number (normalized) |
| BPAMXB | Ø376777777777 | largest positive BPA number |
| EARGMX (EXPMXA) | 88.028 | largest BPA number x such that exp(s) does not overflow |
| EARGMN (EXPMNA) | -88.028 | smallest BPA number such that exp(x) does not underflow |
| LNACC | .0005 | half-length of interval about 1 where LN accuracy decreases |
| LOGACC (LGACC) | .0005 | same for LOG10 |
| SNHACC | 0.5E-10 | same for SINH, except interval is about 0. |
| TNHACC | Ø400400000001 | same for TANH |
| MAXINT | Ø106777777777 | BPA representation of largest FORTRAN integer |
| FRACBD | Ø066575360400 | Smallest positive BPA number such that the low-order digit immediately precedes the radix point. |

A2

APPENDIX B:  PRIMITIVES WRITTEN AT THE

UNIVERSITY OF KANSAS

```
$       GMAP                                                            00000607
        SYMDEF  BPAADD,BPASUB,BPADIV,BPAMUL,BPACEB                      00000608
        BLOCK   BPACOM                                                  00000609
OPTION  BSS     1               THIS IS FOR BROUNDING OPTION            00000610
BPAFLT  BSS     7               THIS WILL RETURN ERROR CODE             00000611
        USE                                                             00000612
 RSSA  BBSS     8                                                       00000613
        ERLK                                                            00000614
SHIFT   BSS     1                                                       00000615
TEST1  EBSS     2                                                       00000616
TEST2   BSS     2                                                       00000617
*                               OKTAB CONTROLS BRANCHING IN BROUND.     00000618
*                               UPPER HALVES OF THE WORDS HOLD BRANCH   00000619
*                               VALUES FOR THE POSITIVE NUMBERS.        00000620
*                               THE LOWER HALVES HOLD JUMP VALUES       00000621
*                               FOR THE NEGATIVE NUMBERS.               00000622
OKTAB   OCT     0               THE FIRST ENTRY IS EMPTY                00000623
        OCT     0               U(+) = 00, U(-) = 00   IN DECIMAL       00000624
        OCT     000030000030    L(+) = 24, L(-) = 24   IN DECIMAL       00000625
        OCT     000030000000    T(+) = 24, T(-) = 00   IN DECIMAL       00000626
        OCT     000013000022    R(+) = 11, R(-) = 18   IN DECIMAL       00000627
        OCT     000000000030    A(+) = 00, A(-) = 24   IN DECIMAL       00000628
*                               EOTAB IS A TABLE WHICH HOLDS THE CORRECT 00000629
*                               FAULTS FOR OVERFLOW BROUNDING.  VALUES  00000630
*                               FOR POSITIVE NUMBERS ARE IN THE UPPER   00000631
*                               HALVES OF THE WORDS, AND THE VALUES FOR 00000632
*                               NEGATIVE NUMBERS ARE IN THE LOWER       00000633
*                               HALVES. 1 IS OVERFLOW. 2 IS INFINITY.   00000634
EOTAB   OCT     0               FIRST ENTRY IS EMPTY                    00000635
        OCT     000002000001    U(+) = INF, U(-) = OV                   00000636
        OCT     000001000002    L(+) =  OV, L(-) = INF                  00000637
        OCT     000001000001    T(+) =  OV, T(-) = OV                   00000638
        OCT     000002000002    R(+) = INF, R(-) = INF                  00000639
        OCT     000002000002    A(+) = INF, A(-) = INF                  00000640
*                               EUTABP AND EUTABN ARE TABLES WHICH      00000641
*                               HOLD THE CORRECT ANSWERS FOR POSITIVE   00000642
*                               AND NEGATIVE UNDERFLOW RESPECTIVELY.    00000643
*                               EU STANDS FOR EXPONENT UNDERFLOW.       00000644
EUTABP  OCT     0               FIRST ENTRY IS EMPTY                    00000645
        OCT     400400000001    U(+) = POSBND                           00000646
        OCT     400000000000    L(+) = ZERO                             00000647
        OCT     400000000000    T(+) = ZERO                             00000648
        OCT     400000000000    R(+) = ZERO EXCEPT FOR EU BY 1          00000649
        OCT     400400000001    A(+) = POSBND                           00000650
EUTABN  OCT     0               FIRST ENTRY IS EMPTY                    00000651
        OCT     400000000000    U(-) = ZERO                             00000652
        OCT     401377777777    L(-) = MAXNEG                           00000653
        OCT     400000000000    T(-) = ZERO                             00000654
        OCT     400000000000    R(-) = ZERO EXCEPT FOR EU BY 1          00000655
        OCT     401377777777    A(-) = MAXNEG                           00000656
ZERO   EOCT     400000000000    ZERO IS 0*(2** -128)                    00000657
        OCT     0               MANTISSA OF DBL PREC IS 0. THIS IS ZERO 00000658
```

```
CHKDBL OCT     000000000377     USED TO CHECK FOR ONES IS DBLE         00000659
       OCT     777777777777     PRECISION PART OF MANTISSA             00000660
MAXPOS OCT     376777777777     EXP IS 127(LARGEST POSS.) MANTISSA     00000661
       OCT     777777777777     IS LARGEST POSITIVE FRACTION.          00000662
MINPOS OCT     400400000000     EXP IS -128, SO THIS IS MIN POS.       00000663
       OCT     0                                                       00000664
POSBND OCT     400400000001     SMALLEST POSITIVE NUMBER ALLOWED       00000665
       OCT     0                                                       00000666
MAXNEG OCT     401377777777     EXP IS -123 AND MANT IS                00000667
       OCT     777777777777     SMALLEST NEG.                          00000668
MINNEG OCT     377000000000     EXP IS 127 AND MAN IS NEG, NORMALIZED  00000669
       OCT     0                MINNEG NOT ALLOWED AS ARGUMENT OF ANSWER 00000670
NEGBND OCT     377000000001     NEGBND IS SINGLE PREC. NEG OF MAXPOS   00000671
       OCT     0                AND IT IS THE MOST NEGATIVE NO. ALLOWED 00000672
UTA    OCT     0                USED FOR U,T AND A OPTIONS             00000673
       OCT     777777777777                                           00000674
POSRND OCT     0                ADD CONST FOR R+                       00000675
       OCT     400000000000                                           00000676
NEGRND OCT     0                ADD CONST FOR R-                       00000677
       OCT     377777777777                                           00000678
EUONE  OCT     376000000000     MAX EXPONENT, UNDERFLOW BY 1           00000679
RSTABL OCT     777,1777,3777,7777,17777,37777,77777,177777,377777     00000680
       OCT     777777,1777777,3777777,7777777,17777777,37777777,7777777700000681
       OCT     177777777,377777777,777777777,1777777777,3777777777    00000682
       OCT     7777777777,17777777777,37777777777,77777777777         00000683
       OCT     177777777777,377777777777,777777777777                 00000684
ADDONEEOCT     0                                                       00000685
   .   OCT     001000000000                                           00000686
SUBONEEOCT     0                                                       00000687
       OCT     1                                                       00000688
BIGCHK OCT     376777777777     USED TO TEST FOR OVERFLOW              00000689
       OCT     0                AGAINST MAXPOS.                        00000690
EUOOFF OCT     747777                                                 00000691
EOON   OCT     020000           EXPONENT OVERFLOW IS BIT 22            00000692
EUON   OCT     010000           EXPONENT UNDERFLOW IS BIT 23           00000693
OFOMSK OCT     000000004000     BIT 24 IS THE OVERFLOW MASK INDICATOR  00000694
TEMP   EBSS    2                                                       00000695
ARG1   EBSS    2                                                       00000696
ARG2   BSS     2                                                       00000697
COUNT  BSS     1                                                       00000698
SAVEXP BSS     1                                                       00000699
SAFE   EBSS    2                                                       00000700
R.U.   BSS     2                                                       00000701
ANS    BSS     2                                                       00000702
BITS   OCT     1000000,2000000,4000000,10000000,20000000,40000000     00000703
       OCT     100000000,200000000,400000000,1000000000,2000000000    00000704
       OCT     4000000000,10000000000,20000000000,40000000000         00000705
       OCT     100000000000,200000000000,400000000000                 00000706
CALL   BSS     1                USED TO IDENTIFY CALLING ROUTINE       00000707
TURN   MACRO                                                           00000708
       STQ     TEMP+1                                                  00000709
       LDQ     #1                                                      00000710
```

B3

```
         STI      TEMP                                                              00000711
         ORSQ     TEMP                                                              00000712
         LDI      TEMP                                                              00000713
         ENDM     TURN                                                              00000714
ENTER    MACRO                                                                      00000715
         SREG     .RSSA                                                             00000716
         STI      .E.L..                                                            00000717
         STX1     .E.L..                                                            00000718
         LDA      OFOMSK    WE WANT TO TURN OFF THE FAULT TRAPS                     00000719
         STI      TEMP      GET A COPY OF THE INDICATORS                            00000720
         ORSA     TEMP      SET OVERFLOW MASK INDICATOR OF COPY                     00000721
         LDA      EUOOFF    TURN EXPONENT OFLO AND UNFLO OFF                        00000722
         ANSA     TEMP                                                              00000723
         LDI      TEMP      FAULT TRAPS ARE NOW DISABLED                           00000724
         LDX2     0,DU      CLEAR ERROR REGS.                                       00000725
         LDX4     0,DU                                                              00000726
         LDX6     0,DU      CLEAR SC                                                00000727
         LDX7     0,DU      CLEAR RI                                                00000728
         EAX0     2,1*      GET FIRST ARG                                           00000729
         IFE      #1,1,21   IF ADD,SUB,MUL,OR DIV, DO NEXT 21 LINES                00000730
         FLD      0,0                                                               00000731
         FCMP     MINNEG    MINNEG IS NOT ALLOWED AS AN ARGUMENT                    00000732
         TNZ      *+2                                                               00000733
         TRA      BROUND    AND GO TO BROUND PART .EO.                             00000734
         FCMP     MINPOS    MINPOS IS NOT ALLOWED AS AN ARG                         00000735
         TNZ      *+3       IF NOT MINPOS, GO ON                                    00000736
         LDX2     3,DU      SET UNDERFLOW FAULT                                     00000737
         TRA      BROUND    GO TO BROUND(.EU.) WITH UNFLO BY 1                     00000738
         FST      ARG1      STORE IN DOUBLE PRECISION AREA                          00000739
         STZ      ARG1+1    CLEAR EXTRA PRECISION                                   00000740
         EAX0     3,1*                                                              00000741
         FLD      0,0                                                               00000742
         FCMP     MINNEG    MINNEG IS NOT ALLOWED AS AN ARGUMENT                    00000743
         TNZ      *+2       GO ON IF ARG. IS NOT MINNEG                            00000744
         TRA      BROUND    AND GO TO BROUND PART .EO.                             00000745
         FCMP     MINPOS    MINPOS IS NOT ALLOWED AS AN ARGUMENT                    00000746
         TNZ      *+3       GO ON IF NOT MINPOS                                     00000747
         LDX2     3,DU      SET UNDERFLOW FAULT                                     00000748
         TRA      BROUND    GO TO BROUND(.EU.) WITH EU BY ONE                     00000749
         FST      ARG2      SAME FOR ARG2                                           00000750
         STZ      ARG2+1    CLEAR EXTRA PRECISION.                                  00000751
         IFE      #1,2,4    IF BPACEB, DO NEXT FOUR LINES                          00000752
         LDA      0,0                                                               00000753
         LDQ      1,0                                                               00000754
         STA      ARG1                                                              00000755
         STQ      ARG1+1                                                            00000756
         LDA      #1,DL     STORE CALL TYPE IN CALL. CEB=2, OTHERS=1               00000757
         STA      CALL                                                             00000758
         ENDM     ENTER                                                            00000759
ARSHFT   MACRO                                                                      00000760
         STX0     TEMP                                                              00000761
         LDX0     COUNT     GET SHIFT COUNT AND SKIP IF ZERO                       00000762
```

B4

```
        TZE        *+14                                                    00000763
        CMPX0      37,DU          SHIFT NORMAL IF .LE. 36 BIT SHIFT         00000764
        TMI        *+11                                                     00000765
        CMPX0      63,DU          SHIFT .GE. 63, CHECK FOR ONES             00000766
        TMI        *+6                                                      00000767
        CANA       =0777777777777 CHECK FOR ANY ONES                        00000768
        TZE        *+2            IF NOT THEN SKIP, ELSE GO ON              00000769
        LDX7       1,DU           SET RESIDUE INDICATOR                     00000770
        DFLD       ZERO           ALL BITS WERE SHIFTED OUT SO LOAD ZERO    00000771
        TRA        *+5            NOW PREPARE TO RETURN                     00000772
        CANA       RSTABL-37,0    CHECK FOR ONES THAT WILL BE LOST          00000773
        TZE        *+2            IF NONE THEN GO SHIFT                     00000774
        LDX7       1,DU           ELSE SET RESIDUE INDICATOR                00000775
        LRL        0,0            SHIFT MANTISSA RIGHT COUNT BITS           00000776
        LDX0       TEMP                                                     00000777
        ENDM       ARSHFT                                                   00000778
NEGU    MACRO                                                              00000779
        DFST       SAFE                                                     00000780
        DFLD       R.U.                                                     00000781
        FNEG                                                                00000782
        DFST       R.U.                                                     00000783
        DFLD       SAFE                                                     00000784
        ENDM       NEGU                                                     00000785
BPACEB  ENTER      2              ENTRY TO CONVERT EXTENDED TO BPA          00000786
        DFLD       ARG1                                                     00000787
        TRA        BROUND                                                   00000788
BPAMUL  ENTER      1                                                        00000789
        FLD        ARG1           MAKE ARGS POSITIVE IF NEC.                00000790
        TPL        BPAM.1         SET SC IF WE CHANGE SIGN                  00000791
        FNEG                                                                00000792
        FST        ARG1                                                     00000793
        LDX6       1,DU                                                     00000794
BPAM.1  FLD        ARG2                                                     00000795
        TPL        BPAM.3                                                   00000796
        FNEG                                                                00000797
        FST        ARG2                                                     00000798
        CMPX6      1,DU           NEED TO FLIP SC.                          00000799
        TZE        BPAM.2                                                   00000800
        LDX6       1,DU                                                     00000801
        TRA        BPAM.3                                                   00000802
BPAM.2  LDX6       0,DU                                                     00000803
BPAM.3  FLD        ARG1           MULTIPLY, DO NOT NORMALIZE                00000804
        UFM        ARG2                                                     00000805
        TEU        BPAMEU                                                   00000806
        TEO        BPAMEO                                                   00000807
        FNO                       EXPONENT IS OK, TRY TO NORMALIZE.         00000808
        TEU        BPAMEU         WE CAN ONLY MAKE THE EXPONENT SMALLER.    00000809
        CMPX6      0,DU           GET SIGN ANSWER CORRECT.                  00000810
        TZE        BROUND                                                   00000811
        FNEG                                                                00000812
        TRA        BROUND                                                   00000813
BPAMEU  FNO                       GET THE ANSWER IN THE BEST FORM POSS.     00000814
```

```
         LXL4    OPTION                                                          00000815
         STZ     BPAFLT                                                          00000816
         CMPX6   0,DU                                                            00000817
         TZE     .EU.          BEFORE TRANSFER, MAKE SURE ANS IS POS.            00000818
         FNEG                                                                    00000819
         TRA     .EU.                                                            00000320
BPAMEO   LXL4    OPTION                                                          00000321
         STZ     BPAFLT                                                          00000822
         DFST    ANS           SAVE THE RESULT AND SEE HOW MUCH WE MUST          00000323
         LDE     0,DU          SHIFT TO NORMALIZE. R(E) WILL HAVE NEG,           00000824
*   NEGATE BEFORE BROUND IF NEEDED                                              00000325
         CMPX6   0,DU                                                            00000826
         TZE     BPAM.4                                                          00000827
         FNEG                                                                    00000828
BPAM.4   FNO                   SHIFT NUM. IT MAY BE ENOUGH TO FIX THE            00000329
         ADE     ARG1          OVERFLOW. THIS FIRST ADD CAN'T OVERFLOW.          00000330
         ADE     ARG2          BECAUSE BOTH EXP. OF VAR. ARE POSITIVE            00000831
         TEO     .EO.                                                            00000832
         TRA     BROUND                                                          00000833
BPADIV   ENTER   1                                                              00000834
         FLD     ARG1          MAKE BOTH MANTISSA POSITIVE                       00000835
         TPL     BPAV.1                                                          00000836
         FNEG                                                                    00000337
         DFST    ARG1                                                            00000838
         LDX6    1,DU                                                            00000839
BPAV.1   FLD     ARG2                                                           00000840
         TZE     ZERODV                                                          00000841
         TPL     BPAV.3                                                          00000842
         FNEG                                                                    00000843
         DFST    ARG2                                                            00000844
         CMPX6   1,DU                                                            00000845
         TZE     BPAV.2                                                          00000846
         LDX6    1,DU                                                            00000847
         TRA     BPAV.3                                                          00000343
BPAV.2   LDX6    0,DU                                                            00000349
BPAV.3   STZ     SHIFT                                                           00000850
         FLD     ARG2          SEE IF DIVIDEND IS LARGER THAN DIVISOR            00000851
         STA     TEMP          THIS WORKS IF BOTH ARGS ARE NORMALIZED            00000852
         FLD     ARG1                                                            00000853
         CMPA    TEMP                                                            00000854
         TMI     BPAV.5                                                          00000855
         TNZ     *+5           IF THE MANTISSAS ARE .5 AND THE SIGN              00000856
         CMPX6   0,DU          WAS CHANGED, THEN WE MUST NOT ADD                 00000857
         TZE     *+3           ONE TO THE SHIFT COUNT                            00000858
         CMPA    =0200000000000 MANTISSA OF ONE HALF                            00000859
         TZE     *+2                                                             00000860
         AOS     SHIFT                                                           00000861
         LRS     1                                                              00000862
         TRA     BPAV.4                                                          00000863
BPAV.5   FLD     ARG1                                                            00000864
BPAV.4   DVF     TEMP                                                            00000365
         CMPQ    0,DU                                                            00000866
```

B6

```
        TZE     SMALL                                                    00000867
        LDX7    1,DU        FILL RESIDUE IND IF THERE IS REMAINDER,      00000868
SMALL   CMPX6   1,DU        SEE IF ANS IS NEG.                           00000869
        TNZ     EXP                                                      00000870
        LDQ     0,DU                                                     00000871
        LDE     0,DU        CLEAR EXPONENT TO PREVENT ERROR IN NEG       00000872
        FNEG                                                             00000873
EXP     LDE     0,DU                                                     00000874
        LDQ     SHIFT                                                    00000875
        CMPQ    0,DU                                                     00000876
        TZE     EXP.1                                                    00000877
        LDQ     0,DU                                                     00000878
        LDE     =0002000,DU                                              00000879
                * BACKSPACE
        LDE     =0002000,DU                                              00003879
EXP.1   FNO                                                              00000880
        DFST    ANS                                                      00000881
        LDA     ARG1                                                     00000882
        ANA     =0776000,DU                                              00000883
        ARS     28                                                       00000884
        LDQ     ARG2                                                     00000885
        ANQ     =0776000,DU                                              00000886
        QRS     28                                                       00000887
        STQ     TEMP                                                     00000888
        SBA     TEMP                                                     00000889
        LDQ     ANS                                                      00000890
        ANQ     =0776000,DU                                              00000891
        QRS     28                                                       00000892
        STQ     TEMP                                                     00000893
        ADA     TEMP                                                     00000894
        CMPA    =128                                                     00000895
        TPL     EXPO        CHECK IF EXP IN RANGE                        00000896
        CMPA    =-128                                                    00000897
        TMI     EXPU                                                     00000898
        ALS     28                                                       00000899
        STA     TEMP                                                     00000900
        DFLD    ANS                                                      00000901
        LDE     TEMP                                                     00000902
        DFST    ANS                                                      00000903
        CMPX7   1,DU        IF THERE WAS A REMAINDER, ADD ONE            00000904
        TNZ     BROUND                                                   00000905
        STE     ADDONE                                                   00000906
        FLD     ADDONE                                                   00000907
        TRA     BROUND                                                   00000908
EXPO    LXL4    OPTION                                                   00000909
        STZ     BPAFLT                                                   00000910
        TRA     .EO.                                                     00000911
EXPU    LXL4    OPTION                                                   00000912
        STZ     BPAFLT                                                   00000913
        TRA     .EU.                                                     00000914
ZERODV  LDA     4,DL        DIVISION BY ZERO CASE                        00000915
        STA     BPAFLT                                                   00000916
        TRA     EXIT                                                     00000917
BPAADD  ENTER   1                                                        00000918
```

```
        DFLD    ARG2            GET ARG2, TEST FOR ZERO, STORE IN R.U.    00000919
        INZ     *+3                                                      00000920
        FLD     ARG1            ARG2 = 0, LOAD ARG1 AND RETURN            00000921
        TRA     EXIT                                                     00000922
        DFST    R.U.                                                     00000923
SUBIN   DFLD    ARG1            BPASUB COMES HERE. GET ARG1.             00000924
        TNZ     ADD1            IF ARG1 = 0, PREPARE TO LEAVE BPAADD      00000925
        DFLD    R.U.            OUR ANSWER IS ARG2.                      00000926
        TRA     EXIT            GO TO EXIT, ONE ARG WAS ZERO             00000927
ADD1    DFST    TEST1                                                    00000928
        DFLD    R.U.            WE WILL MAKE BOTH POSITIVE AND COMPARE    00000929
        TPL     ADD2            DFCMG DOES NOT WORK (NOV.1976)            00000930
        FNEG                    MAKE 2ND ARG POSITIVE FOR DFCMP          00000931
ADD2    DFST    TEST2           NOW RELOAD TEST1 AND MAKE +              00000932
        DFLD    TEST1                                                    00000933
        TPL     ADD3                                                     00000934
        FNEG                    MAKE R(EAQ) +                            00000935
ADD3    DFCMP   TEST2           MAKE THE COMPARISON                      00000936
        TZE     ADD5            IF EQUAL, DO NOT SWITCH                   00000937
        TMI     ADD5            IF R(EAQ) .LT. R.U., LEAVE               00000938
ADD4    DFLD    ARG1            R(EAQ) MAY BE WRONG SIGN, SO RELOAD ARG1 00000939
        DFST    R.U.            AND SWITCH R(EAQ) AND R.U.               00000940
        DFLD    ARG2                                                     00000941
        TRA     ADD6                                                     00000942
ADD5    DFLD    ARG1            RECONSTITUTE THE SIGN OF R(EAQ)          00000943
ADD6    DFCMP   ZERO            R(EAQ) MUST BE SET POSITIVE              00000944
        TPL     ADD7                                                     00000945
        FNEG                    NEGATE R(EAQ) AND R.U., SET SIGN CHANGE  00000946
        NEGU                                                            00000947
        LDX6    1,DU                                                     00000948
ADD7    STE     SAVEXP          GET EXPONENTS, COMPUTE SHIFT COUNT       00000949
        LDX0    R.U.                                                     00000950
        ANX0    =0776000,DU                                              00000951
        SBX0    SAVEXP                                                   00000952
        TOV     *+2             IF OVERFLOW, THE SHIFT NUMBER            00000953
        TRA     *+2             WAS TOO BIG, AND SO WE LOAD 64           00000954
        LDX0    =0200000,DU     AS THE SHIFT COUNT                       00000955
        TMI     ADD8            IF SHIFT IS NEGATIVE NO SHIFT NEEDED     00000956
        TZE     ADD8            NO SHIFT NEEDED, SHIFT COUNT WAS ZERO    00000957
        EAQ     0,0                                                     00000958
        QRL     10                                                      00000959
        STQ     COUNT           STORE SHIFT COUNT IN COUNT              00000960
        LDQ     =0                                                      00000961
        ARSHFT  COUNT           SHIFT MANTISSA RIGHT COUNT BITS          00000962
        LDE     R.U.            LOAD CORRECT EXPONENT FOR R(EAQ)         00000963
ADD8    DFAD    R.U.            ADD THE TWO NUMBERS AND NORMALIZE.       00000964
        CMPX6   0,DU            TEST FOR SIGN CORRECTION                 00000965
        TZE     ADD11           AND SKIP TO ADD11 IF NONE NEEDED         00000966
        TEO     *+3             GO TO CORRECTION FOR EXPONENT OFLO       00000967
        TEU     ADD9            SAME FOR EXPONENT UNDERFLOW              00000968
        TRA     ADD10           AND AGAIN IF NO FAULT OCCURRED           00000969
*                               MINPOS IS THE RESULT OF ONE OVERFLOW     00000970
```

```
*                              CORRECT NEGATIVE IS MINNEG, BUT NEGATIVE 00000971
*                              MINPOS GIVES UNDERFLOW. MINNEG IS NOT    00000972
*                              ALLOWED IN INTERVAL, SO GIVE IT OVERFLOW 00000973
*                              TEU TURN OFF UNDERFLOW IF WE NEGATED     00000974
*                              MINPOS.                                  00000975
        FNEG                   NEGATE THE MANTISSA                      00000976
        TEU      *+1                                                    00000977
        TURN     EOON          TURN EO BACK ON                         00000978
        TRA      ADD11         WE ARE THROUGH NEGATING OVERFLOW         00000979
*                              MINNEG IS THE RESULT OF ONE UNDERFLOW;   00000980
*                              CORRECT NEGATIVE IS MINPOS, BUT NEGATIVE 00000981
*                              MINNEG GIVES OVERFLOW.  MINPOS IS NOT    00000982
*                              ALLOWED IN INTERVAL, SO GIVE IT UNFLOW.  00000983
*                              TEO TURNS THE OVERFLOW INDICATOR OFF IF  00000984
*                              WE NEGATED MINNEG.                       00000985
ADD9    FNEG                   NEGATE THE UNDERFLOW ANSWER              00000986
        TEO      *+1           TURN OFLO OFF IF IT IS ON                00000987
        TURN     EUON          TURN EU INDICATOR BACK ON                00000988
        TRA      ADD11         AND GO TO ADD11                          00000989
ADD10   FNEG                   NO PROBLEM, JUST NEGATE                  00000990
ADD11   CMPX7    0,DU          CHECK RESIDUE INDICATOR, EXIT IF ZERO    00000991
        TNZ      ADD12                                                  00000992
        DFCMP    ZERO          EXIT IF ZERO, ELSE GO TO BROUND          00000993
        TZE      EXIT                                                   00000994
        TRA      BROUND                                                 00000995
*                                                                      00000996
*                                                                      00000997
*       AT THIS POINT WE MAY NEED A RESIDUE CORRECTION.                00000998
*       WE HAVE THE FOUR FOLLOWING POSSIBLE CASES..                    00000999
*                1. SC = 0 AND R(EAQ) .GE. 0                           00001000
*                2. SC = 0 AND R(EAQ) .LT. 0                           00001001
*                3. SC = 1 AND R(EAQ) .GE. 0                           00001002
*                4. SC = 1 AND R(EAQ) .LT. 0                           00001003
*                                                                      00001004
*       CASES 1 AND 2.  IF THE EXTRA PRECISION WORD                    00001005
*                       IS 0, SET IT TO 1,                             00001006
*                       AND THEN GO TO BROUND.                         00001007
*       CASES 3 AND 4.  IF THE EXTRA PRECISION WORD                    00001008
*                       IS 0, SUBTRACT 1 FROM THE                      00001009
*                       FULL MANTISSA, AND THEN GO                     00001010
*                       TO BROUND.                                     00001011
*                                                                      00001012
*                                                                      00001013
ADD12   CMPX6    0,DU          IDENTIFY THE CASE. CHECK SC             00001014
        TNZ      ADD13         IF SC = 1, GO TO ADD13                  00001015
        DFST     TEMP                                                  00001016
        LDQ      TEMP+1                                                00001017
        TNZ      BROUND        IF Q NOT ZERO, WE ARE DONE HERE         00001018
        LDQ      =1                                                    00001019
        STQ      TEMP+1        SET Q TO 1 REPLACE IN TEMP              00001020
        DFLD     TEMP          LOAD TEMP AND GO TO BROUND              00001021
        TRA      BROUND                                                00001022
```

```
AOD13    DFST     TEMP          WE HAVE CASE 3 OR 4               00001023
         LDQ      TEMP+1        CHECK FOR NONZERO AND LEAVE IF SO 00001024
         TNZ      BROUND                                         00001025
         DFLD     TEMP          ELSE SUBTRACT ONE FROM AQ        00001026
         STE      SUBONE                                         00001027
         DFSB     SUBONE                                         00001028
         TRA      BROUND                                         00001029
BPASUB   ENTER    1                                             00001030
         DFLD     ARG2          LOAD ARG2, TEST FOR ZERO, NEGATE 00001031
         TNZ      *+3                                            00001032
         FLD      ARG1                                          00001033
         TRA      EXIT                                          00001034
         FNEG                   NEGATE AND GO TO  SUBIN IN BPAADD 00001035
         DFST     R.U.                                          00001036
         DFST     ARG2                                          00001037
         TRA      SUBIN                                         00001038
BROUND   FNO                    NORMALIZE THE R(EAQ) IF IT IS NOT 00001039
         STZ      BPAFLT        PREPARE RETURN VARIABLE          00001040
         DFCMP    ZERO          IF ANSWER IS ZERO WE ARE DONE    00001041
         TZE      EXIT          THEN GO TO EXIT,WE ARE DONE      00001042
         LXL4     OPTION                                         00001043
         TMI      *+4                                            00001044
         TZE      *+3                                            00001045
         CMPX4    6,DU                                          00001046
         TMI      *+2                                            00001047
         LDX4     4,DU                                          00001048
         TEO      .EO.          DID OP CAUSE EXP OVERFLOW OR UNDERF. 00001049
         TEU      .EU.                                          00001050
         DFCMP    MINPOS        IF ANSWER IS MINPOS, GO TO .EU.1 00001051
         TZE      .EU.1                                         00001052
         DFCMP    NEGBND        IF ANS.>=NEGBND,NO ERROR         00001053
         TMI      .EO.                                          00001054
         DFCMP    BIGCHK        IF BIGGER,THEN OFLOW EVEN THOUGH 00001055
         TMI      *+3           IT IS SMALLER THAN MAXPOS        00001056
         TZE      BRND.1        GO LOAD OKTAB JUMP VALUE.        00001057
         TRA      .EO.          WE HAVE OVERFLOW, WITHOUT EO ON. 00001058
         CANAQ    CHKDBL        IF NO ONES IN DBL PREC PART OF MANT 00001059
         TZE      EXIT          THEN WE NEED NO BROUNDING        00001060
         DFCMP    ZERO          CHOOSE UPPER OR LOWER HALF OF WORD 00001061
         TMI      *+3           TAKE LOWER HALF IF MINUS         00001062
BRND.1   LDX4     OKTAB,4       GET JUMP VALUE FROM TABLE        00001063
         TRA      NOFLT,4       GO TO RIGHT BROUNDING            00001064
         LXL4     OKTAB,4       GET JUMP VALUE FROM TABLE        00001065
         TRA      NOFLT,4       GO TO BROUNDING FOR MINUS        00001066
.EO.     FCMP     ZERO          TEST WHETHER POSITIVE OR NEGATIVE 00001067
         TMI      .EO.1         IF NEGATIVE GO TO .EO.1          00001068
         LDX2     EOTAB,4       LOAD FAULT VECTOR FROM TABLE     00001069
         FLD      MAXPOS        LOAD POSITIVE ANSWER             00001070
         TRA      EXIT                                          00001071
.EO.1    FLD      NEGBND        LOAD NEGATIVE ANSWER             00001072
         LXL2     EOTAB,4       LOAD FAULT VECTOR               00001073
         TRA      EXIT          WE ARE THROUGH                  00001074
```

B10

```
.U       LDX2     3,DU              LOAD UNDERFLOW FOR ALL OPTIONS         00001075
         FCMP     ZERO              TEST IF POSITIVE OR NEGATIVE           00001076
         TMI      .EU.2             IF NEGATIVE, GO TO .EU.2               00001077
         CMPX4    4,DU              ROUNDING TAKES SPECIAL TREATMENT       00001078
         TZE      *+3               SO SKIP THESE TWO LINES                00001079
         FLD      EUTABP,4          LOAD CORRECT ANSWER FOR U,L,T,A        00001080
         TRA      EXIT              AND GO TO EXIT                         00001081
         STE      SAVEXP            CHECK IF UNFLO BY 1                     00001082
         LDX0     SAVEXP            R(EAQ) NOW HAS ZERO FOR R OPTION       00001083
         CMPX0    EUONE             IF UNFLO BY 1, LOAD POSBND             00001084
         TNZ      *+3                                                      00001085
.EU.1    FLD      POSBND                                                   00001086
         TRA      EXIT              WE ARE DONE WITH POSITIVE CASE         00001087
         FLD      ZERO              IF EU BY >1 LOAD ZERO                  00001088
         TRA      EXIT              AND GO TO EXIT                         00001089
.EU.2    CMPX4    4,DU              IF NOT ROUND, LOAD VALUE FROM TABLE    00001090
         TZE      *+3                                                      00001091
         FLD      EUTABN,4          LOAD RIGHT ANSWER FOR U,L,T,A          00001092
         TRA      EXIT              AND GO TO EXIT                         00001093
         STE      SAVEXP            CHECK IF UNFLO BY 1                     00001094
         LDX0     SAVEXP            R(EAQ) NOW HAS ZERO FOR R OPTION       00001095
         CMPX0    EUONE             IF UNFLO BY 1, LOAD MAXNEG             00001096
         TNZ      *+3               ELSE GO LOAD ZERO                      00001097
         FLD      MAXNEG                                                   00001098
         TRA      EXIT                                                     00001099
         FLD      ZERO              LOAD ZERO AS THE ANSWER, EU > 1        00001100
         TRA      EXIT              AND GO TO EXIT                         00001101
*                                  THE FIRST BROUNDING WE DO IS U(+,-),   00001102
*                                  T(-), AND A(+).  IN ALL CASES, WE      00001103
*                                  ADD UTA, WHICH GENERATES A CARRY       00001104
*                                  INTO THE SINGLE PREC. ANSWER IF THERE  00001105
*                                  ARE ANY 1'S IN THE EXTRA PREC.WORD     00001106
NOFLT    STE      UTA                                                      00001107
         DFAD     UTA                                                      00001108
         TEO      *+3               OFLO IS POSITIVE. HAVE U OR A          00001109
         TEU      *+5               UNDERFLOW IS NEGATIVE, T               00001110
         TRA      EXIT              GO TO EXIT IF ANSWER OK                00001111
         LDX2     2,DU              SET INFINITY FOR OVERFLOW              00001112
         DFLD     MAXPOS            WANT MAXPOS AS ANSWER                  00001113
         TRA      EXIT                                                     00001114
         LDX2     3,DU              SET UNDERFLOW                          00001115
         DFLD     ZERO              LOAD ZERO ON T OPTION                  00001116
         TRA      EXIT              END OF UTA BROUNDING                   00001117
*                                  ROUND POSITIVE NUMBER. BIT 28 MUST     00001118
*                                  BE 1 TO GO UP, ELSE CHOP OFF EXTRA     00001119
*                                  PRECISION.                             00001120
RPOS     STE      POSRND                                                   00001121
         DFAD     POSRND                                                   00001122
         TEO      *+2               OVERFLOW IS POSITIVE                   00001123
         TRA      EXIT                                                     00001124
         LDX2     2,DU              SET INFINITY                          00001125
         DFLD     MAXPOS            WANT MAXPOS AS ANSWER                  00001126
```

B11

```
        TRA     EXIT                                                    00001127
*                               ROUND A NEGATIVE ANSWER. ROUND          00001128
*                               VALUE IN MIDDLE TOWARD ZERO.            00001129
*                               BIT 26 MUST BE 0 TO GO AWAY            00001130
*                               FROM ZERO.                             00001131
RNEG    STE     NEGRND                                                  00001132
        DFAD    NEGRND                                                  00001133
        TEU     *+2             UNDERFLOW IS NEGATIVE                   00001134
        TRA     EXIT                                                    00001135
        LDX2    3,DU            SET UNDERFLOW, UNFLO BY 1               00001136
        DFLD    MAXNEG          SO LOAD MAXNEG, NOT 0                   00001137
EXIT    SXL2    BPAFLT          RETURN FAULT VECTOR IN COMMON           00001138
        LXL2    CALL                                                    00001139
        CMPX2   2,DU                                                    00001140
        TZE     EXIT.1                                                  00001141
        EAX0    4,1*                                                    00001142
        FST     0,0             AND BROUNDED ANSWER IN ARGUMENT         00001143
        TRA     EXIT.2                                                  00001144
EXIT.1  EAX0    3,1*                                                    00001145
        FST     0,0                                                     00001146
EXIT.2  LREG    .RSSA                                                   00001147
        RET     .E.L..                                                  00001148
        END                                                             00001149
```

In accordance with letter from DAEN-RDC, DAEN-ASI dated
22 July 1977, Subject:  Facsimile Catalog Cards for
Laboratory Technical Publications, a facsimile catalog
card in Library of Congress MARC format is reproduced
below.